

TESTING AND NON-LINEAR PRECONDITIONING OF THE PROXIMAL POINT METHOD

Tuomo Valkonen*

2017-03-16

Abstract Employing the ideas of non-linear preconditioning and testing of the classical proximal point method, we formalise common arguments in convergence rate and convergence proofs of optimisation methods to the verification of a simple iteration-wise inequality. When applied to fixed point operators, the latter can be seen as a generalisation of firm non-expansivity or the α -averaged property. The main purpose of this work is to provide the abstract background theory for our companion paper “Block-proximal methods with spatially adapted acceleration”. In the present account we demonstrate the effectiveness of the general approach on several classical algorithms, as well as their stochastic variants. Besides, of course, the proximal point method, these method include the gradient descent, forward–backward splitting, Douglas–Rachford splitting, Newton’s method, as well as several methods for saddle-point problems, such as the Alternating Directions Method of Multipliers, and the Chambolle–Pock method.

Get the version from <http://tuomov.iki.fi/publications>, citations broken in this one due broken arXiv biblatex support.

1 INTRODUCTION

The proximal point method for monotone operators [17, 22], while infrequently used by itself, can be found as a building block of many popular optimisation algorithms. Indeed, many important application problems can be written in the form

$$(P) \quad \min_x G(x) + F(Kx)$$

for convex non-smooth G and F , and a linear operator K . Examples abound in image processing and data science. The problem (P) can often be solved by methods such as forward–backward splitting, ADMM (alternating directions method of multipliers) and their variants [2, 16, 11, 6]. They all involve a proximal point step.

The equivalent saddle point form of (P) is

$$(S) \quad \min_x \max_y G(x) + \langle Kx, y \rangle - F^*(y).$$

In particular within mathematical image processing and computer vision, a popular algorithm for solving (S) is the primal–dual method of Chambolle and Pock [6]. As discovered in [12], the

*Department of Mathematical Sciences, University of Liverpool, United Kingdom. tuomo.valkonen@iki.fi

method can most concisely be written as a *preconditioned proximal point method*, solving on each iteration for $u^{i+1} = (x^{i+1}, y^{i+1})$ the variational inclusion

$$(PP_0) \quad 0 \in H(u^{i+1}) + M_{i+1}(u^{i+1} - u^i),$$

where the monotone operator

$$H(u) := \begin{pmatrix} \partial G(x) + K^*y \\ \partial F^*(y) - Kx \end{pmatrix}$$

encodes the optimality condition $0 \in H(\widehat{u})$ for (S). In the standard proximal point method [22], one would take $M_{i+1} = I$ the identity. With this choice, (PP₀) is generally difficult to solve. In the Chambolle–Pock method the *preconditioning operator* is given for suitable step length parameters $\tau_i, \sigma_{i+1}, \theta_i > 0$ by

$$(1.1) \quad M_{i+1} := \begin{pmatrix} \tau_i^{-1}I & -K^* \\ -\theta_i K & \sigma_{i+1}^{-1}I \end{pmatrix}.$$

This choice of M_{i+1} decouples the primal x and dual y updates, making the solution of (PP₀) feasible in a wide range of problems. If G is strongly convex, the step length parameters $\tau_i, \sigma_{i+1}, \theta_i$ can be chosen to yield $O(1/N^2)$ convergence rates of an ergodic duality gap and the quadratic distance $\|x^i - \widehat{x}\|^2$.

In our earlier work [25], we have modified M_{i+1} as well as the condition (PP₀) to still allow a level of mixed-rate acceleration when G is strongly convex only on sub-spaces. Our convergence proofs were based on *testing* the abstract proximal point method by a suitable operator, which encodes the desired and achievable convergence rates on relevant subspaces.

In the present paper, we extend this theoretical approach to non-linear preconditioning, non-invertible step-length operators, and arbitrary monotone operators H . Our main purpose is to provide the abstract background theory for our companion paper [24]. Here, within these pages, we demonstrate that several classical optimisation methods—including the second-order Newton’s method—can also be seen as variants of the proximal point method, and that their common convergence rate and convergence proofs reduce to the verification of a simple iteration-wise inequality. Through application of our theory to Browder’s fixed point theorem [4] in Section 2.5, we see that our inequality generalises the concepts of firm non-expansivity or the α -averaged property. Our theory also covers stochastic variants of the considered algorithms.

In Section 2, we start by developing our theory for general monotone operators H . This extends, simplifies, and clarifies the more disconnected results from [25] that concentrated on saddle-point problems with preconditioners derived from (1.1). We demonstrate our results on the basic proximal point method, gradient descent, forward–backward splitting, Douglas–Rachford splitting, and Newton’s method. The proximal step in forward–backward splitting and proximal Newton’s method can be introduced completely “free”, without any additional proof effort, in our approach.

In Section 3 we specialise our work to saddle-point problems, and demonstrate the results on variants of the Chambolle–Pock method, ADMM, and the Generalised Iterative Soft Thresholding (GIST) algorithm of [16]. In the final Section 4 we extend our results and examples to produce the convergence of ergodic duality gaps. This is also where we move to the stochastic

setting, which allows our results to be used to study various stochastic block-coordinate descent methods. We refer to [26] for a review of this class of methods. In the companion paper [24], we will apply our results to stochastic primal-dual methods with coordinate-wise adapted step lengths.

Besides already cited works, other previous work related to ours includes that on generalised proximal point methods, such as [5, 8], as well inertial methods for variational inclusions [15].

2 AN ABSTRACT PRECONDITIONED PROXIMAL POINT ITERATION

2.1 NOTATION AND GENERAL SETUP

We use $C(X)$ to denote the space of convex, proper, lower semicontinuous functions from X to the extended reals $\overline{\mathbb{R}} := [-\infty, \infty]$, and $\mathcal{L}(X; Y)$ to denote the space of bounded linear operators between Hilbert spaces X and Y . We denote the identity operator by I . For $T, S \in \mathcal{L}(X; X)$, we write $T \geq S$ when $T - S$ is positive semidefinite. Also for possibly non-self-adjoint T , we introduce the inner product and norm-like notations

$$(2.1) \quad \langle x, z \rangle_T := \langle Tx, z \rangle, \quad \text{and} \quad \|x\|_T := \sqrt{\langle x, x \rangle_T}.$$

For a set $A \subset \mathbb{R}$, we write $A \geq 0$ if every element $t \in A$ satisfies $t \geq 0$.

Our overall wish is to find some $\widehat{u} \in U$, on a Hilbert space U , solving for a given set-valued map $H : U \rightrightarrows U$ the variational inclusion

$$(2.2) \quad 0 \in H(\widehat{u}).$$

In the present [Section 2](#), H will be arbitrary, but in [Section 3](#), where we specialise the results, and in [Section 4](#), where we consider gap estimates, we concentrate on H arising from the saddle point problem (S).

Our strategy towards finding a solution \widehat{u} is to introduce an arbitrary non-linear iteration-dependent *preconditioner* $V_{i+1} : U \rightarrow U$ and a *step length operator* $W_{i+1} \in \mathcal{L}(U; U)$. With these, we define the generalised proximal point method, which on each iteration $i \in \mathbb{N}$ solves for u^{i+1} from

$$(PP) \quad 0 \in W_{i+1}H(u^{i+1}) + V_{i+1}(u^{i+1}),$$

We assume that V_{i+1} splits into $M_{i+1} \in \mathcal{L}(U; U)$, and $V'_{i+1} : U \rightarrow U$ as

$$(2.3) \quad V_{i+1}(u) = V'_{i+1}(u) + M_{i+1}(u - u^i).$$

More generally, to rigorously extend our approach to cases that would otherwise involve set-valued V_{i+1} , we also consider for $\widetilde{H}_{i+1} : U \rightrightarrows U$ the iteration

$$(PP^\sim) \quad 0 \in \widetilde{H}_{i+1}(u^{i+1}) + M_{i+1}(u^{i+1} - u^i).$$

2.2 BASIC ESTIMATES

We analyse (PP) and (PP[~]) by applying a *testing operator* $Z_{i+1} \in \mathcal{L}(U; U)$, following the ideas introduced in [25]. The product $Z_{i+1}M_{i+1}$ with the linear part of the preconditioner, will, as we soon demonstrate, be an indicator of convergence rates.

Theorem 2.1. *On a Hilbert space U , let $\tilde{H}_{i+1} : U \rightrightarrows U$, and $M_{i+1}, Z_{i+1} \in \mathcal{L}(U; U)$ for $i \in \mathbb{N}$. Suppose (PP[~]) is solvable, and denote the iterates by $\{u^i\}_{i \in \mathbb{N}}$. If $Z_{i+1}M_{i+1}$ is self-adjoint, and*

$$(CI^{\sim}) \quad \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}}^2 \\ + \langle \tilde{H}_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}} \geq -\Delta_{i+1}(\widehat{u})$$

for all $i \in \mathbb{N}$ and some $\widehat{u} \in U$, then

$$(DI) \quad \frac{1}{2} \|u^N - \widehat{u}\|_{Z_{N+1}M_{N+1}}^2 \leq \frac{1}{2} \|u^0 - \widehat{u}\|_{Z_1M_1}^2 + \sum_{i=0}^{N-1} \Delta_{i+1}(\widehat{u}) \quad (N \geq 1).$$

Corollary 2.2. *On a Hilbert space U , let $H : U \rightrightarrows U$. Also let $Z_{i+1}, W_{i+1}, M_{i+1} \in \mathcal{L}(U; U)$, and $V'_{i+1} : U \rightarrow U$ for $i \in \mathbb{N}$. Suppose (PP) is solvable for V_{i+1} as in (2.3). Denote the iterates by $\{u^i\}_{i \in \mathbb{N}}$. Let $\widehat{u} \in H^{-1}(0)$. If $Z_{i+1}M_{i+1}$ is self-adjoint, and*

$$(CI) \quad \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}}^2 \\ + \langle W_{i+1}[H(u^{i+1}) - H(\widehat{u})] + V'_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}} \geq -\Delta_{i+1}(\widehat{u}),$$

for all $i \in \mathbb{N}$, then (CI[~]) and (DI) hold for $\tilde{H}_{i+1}(u) := W_{i+1}H(u) + V'_{i+1}(u)$.

For (PP), the condition (CI) is often more practical to verify than (CI[~]) thanks to the additional structure introduced by $H(\widehat{u}) \ni 0$. Indeed, in many of our examples, we can eliminate H through monotonicity. To derive gap estimates in Section 4, we will however need (CI[~]).

Proof of Theorem 2.1. Inserting (PP[~]) into (CI[~]), we obtain

$$(2.4) \quad \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}}^2 \\ - \langle u^{i+1} - u^i, u^{i+1} - \widehat{u} \rangle_{Z_{i+1}M_{i+1}} \geq -\Delta_{i+1}(\widehat{u}).$$

We recall for general self-adjoint M the three-point formula

$$(2.5) \quad \langle u^{i+1} - u^i, u^{i+1} - \widehat{u} \rangle_M = \frac{1}{2} \|u^{i+1} - u^i\|_M^2 - \frac{1}{2} \|u^i - \widehat{u}\|_M^2 + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_M^2.$$

Using this with $M = Z_{i+1}M_{i+1}$, we rewrite (2.4) as

$$\frac{1}{2} \|u^i - \widehat{u}\|_{Z_{i+1}M_{i+1}}^2 - \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+2}M_{i+2}}^2 \geq -\Delta_{i+1}(\widehat{u}).$$

Summing this over $i = 0, \dots, N-1$, we obtain (DI). \square

Remark 2.3 (Bregman distances). *The three-point formula (2.5) generalises to Bregman distances [8]. If $Z_{i+1} = \phi_{i+1}I$ for some scalar ϕ_{i+1} , it is then easy to generalise [Theorem 2.1](#) from $\frac{1}{2}\|\cdot\|_{M_{i+1}}$ to more general Bregman distances. While we do occasionally work with V_{i+1} arising as the gradient of a more general Bregman distance, we will, however, not benefit from more general M_{i+1} .*

The next two results demonstrate how the estimate of [Theorem 2.1](#) can be used to prove convergence with or without rates.

Proposition 2.4 (Convergence with a rate). *Suppose (DI) holds with $\Delta_{i+1}(\widehat{u}) \leq 0$, and that $Z_{N+1}M_{N+1} \geq \mu(N)I$. Then $\|u^N - \widehat{u}\|^2 \rightarrow 0$ at the rate $O(1/\mu(N))$.*

Proof. Immediate from (DI). □

Proposition 2.5 (Weak convergence). *Suppose $Z_i M_i = Z_0 M_0 \geq 0$ is self-adjoint, and that the iterates of (PP \sim) satisfy (CI \sim) with $\Delta_{i+1}(\widehat{u}) \leq -\frac{\delta}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$ for all $\widehat{u} \in \widehat{U} := \{u \in U \mid 0 \in H(u)\}$ and some $\delta > 0$. If*

$$(CL) \quad Z_{i+1}M_{i+1}(u^{i+1} - u^i) \rightarrow 0 \text{ and } u^{i^k} \rightharpoonup u \implies \limsup_{k \rightarrow \infty} \widetilde{H}_{i+1}(u^{i^k}) \subset W_*H(u)$$

for some non-singular $W_* \in \mathcal{L}(U; U)$, then $Z_0 M_0(u^i - u^*) \rightarrow 0$ weakly in U for some u^* satisfying $0 \in H(u^*)$.

The \limsup denotes the (strong) outer limit [?, see, e.g.,]rockafellar-wets-va. For the proof, we use the next lemma. Its earliest version is contained in the proof of [18, Theorem 1].

Lemma 2.6 ([4, Lemma 6]). *On a Hilbert space X , let $\widehat{X} \subset X$ be closed and convex, and $\{x^i\}_{i \in \mathbb{N}} \subset X$. If the following conditions hold, then $x^i \rightharpoonup x^*$ weakly in X for some $x^* \in \widehat{X}$:*

- (i) $i \mapsto \|x^i - x^*\|$ is non-increasing for all $x^* \in \widehat{X}$.
- (ii) All weak limit points of $\{x^i\}_{i \in \mathbb{N}}$ belong to \widehat{X} .

Proof of Proposition 2.5. Since $Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2} \leq 0$, it is easy to see that (CI \sim) and consequently (DI) holds for all $\widehat{u} \in U' := \text{cl conv } \widehat{U}$. We apply [Theorem 2.1](#) on any $\widehat{u} \in U'$. Using $\Delta_{i+1}(\widehat{u}) \leq -\frac{\delta}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$, we have $Z_{i+1}M_{i+1}(u^{i+1} - u^i) \rightarrow 0$. By (PP \sim) and (CL), any weak limit point u^* of the sequence $\{u^i\}_{i \in \mathbb{N}}$ then satisfies $u^* \in \widehat{U} \subset U'$. Since $A := Z_0 M_0 = Z_{i+1}M_{i+1}$, this verifies condition (ii) of the lemma for $x^i := A^{1/2}u^i$ and $X' := A^{1/2}\widehat{U}$ on $X := A^{1/2}U \subset U$. Applied with $N = 1$ and u^i in place of u^0 , (DI) shows condition (i) of the lemma. Thus $x^i \rightharpoonup x^* \in \widehat{X}$. But $x^* = A^{1/2}u^*$ for some $u^* \in \widehat{U}$. Thus $A(u^i - u^*) \rightarrow 0$. This implies $Z_0 M_0(u^i - u^*) \rightarrow 0$ weakly. □

2.3 EXAMPLES OF FIRST-ORDER METHODS

We now look at several concrete examples.

Example 2.1 (The proximal point method). Take $M_i = I$, $V'_i = 0$, and $W_{i+1} = \tau_i I$ for some $\tau_i > 0$. Then (PP) is the standard proximal point method with step length $1/\tau_i$. If H is

maximal monotone, $\{u^i\}_{i \in \mathbb{N}}$ converges weakly to some $u^* \in H^{-1}(0)$.

Proof of convergence. We take $Z_{i+1} = \phi_i I$ for some $\phi_i > 0$. As long as $\phi_i \geq \phi_{i+1}$, the monotonicity of H clearly shows (CI) with $\Delta_{i+1}(\widehat{u}) = -\frac{\phi_i}{2} \|u^{i+1} - u^i\|^2$. Using the maximal monotonicity, Minty's theorem [?, e.g.,] Theorem 21.1]bauschke2011convex guarantees the solvability of (PP). Thus the conditions of Corollary 2.2 are satisfied. Maximal monotonicity also guarantees that H is weak-to-strong outer semicontinuous; see Lemma A.1. This establishes (CL). Taking $\phi_i \equiv \phi_0$ for constant $\phi_0 > 0$, so that $Z_{i+1}M_{i+1} = Z_0M_0 = \phi_0 I$, it remains to refer to Proposition 2.5. \square

Example 2.2 (Accelerated proximal point method). Continuing from Example 2.1, suppose H is strongly monotone. Then $\langle H(u^{i+1}) - H(\widehat{u}), u^{i+1} - \widehat{u} \rangle \geq \gamma \|u^{i+1} - \widehat{u}\|^2$ for some $\gamma > 0$, so (CI) continues to hold with $\Delta_{i+1}(\widehat{u}) = -\frac{\phi_i}{2} \|u^{i+1} - u^i\|^2$ if $\phi_i(1 + 2\gamma\tau_i) \geq \phi_{i+1}$. This is the case for $\tau_{i+1} := \tau_i/\sqrt{1 + 2\gamma\tau_i}$, and $\phi_{i+1} := 1/\tau_{i+1}^2$. The testing variable ϕ_N is of the order $\Theta(N^2)$ [6, 25], so we get convergence of $\|u^N - \widehat{u}\|^2$ to zero at the rate $O(1/N^2)$ from Corollary 2.2 and Proposition 2.4.

To facilitate the analysis algorithms with a proximal step, we introduce the following strengthened version of (CI), assumed to hold for some $\Delta_{i+1}(u^*; u)$ at all $u \in U_{i+1} \subset U$ and $u^* \in U$:

$$(CI^*) \quad \frac{1}{2} \|u - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2} \|u - u^*\|_{Z_{i+1}M_{i+1} - Z_{i+2}M_{i+2}}^2 + \langle W_{i+1}(H(u) - H(u^*)) + V'_{i+1}(u), u - u^* \rangle_{Z_{i+1}} \geq -\Delta_{i+1}(u^*; u).$$

Note that only the choice $u = u^{i+1}$ and $u^* = \widehat{u}$ implies (CI \sim) and thus convergence. The role of the subset U_{i+1} is to model a compatible range of u^{i+1} between $H = A$ and $H = A + B$ in the next lemma. Typically $U_{i+1} = U$, but for the stochastic examples of Section 4.5, we will need to make restrictions.

Lemma 2.7. *Let $A, B : U \rightrightarrows U$. Suppose (CI \sim) holds for $H = A$, and that*

$$(2.6) \quad \langle B(u) - B(u^*), u - u^* \rangle_{Z_{i+1}W_{i+1}} \geq 0, \quad (u \in U_{i+1}, u^* \in U).$$

Then (CI \sim) holds for $H = A + B$ with W_{i+1} , M_{i+1} , Z_{i+1} , V'_{i+1} and $\Delta_{i+1}(u, u^)$ unchanged. Moreover, if v^{i+1} solves (PP) for $H = A$, then $u^{i+1} := (I + W_{i+1}B)^{-1}(v^{i+1})$ solves (PP) for $H = A + B$.*

Proof. Using (2.6), B is easily eliminated from (CI \sim). The result is (CI \sim) for $H = A$. The relationship between v^{i+1} and u^{i+1} is immediate from expansion of (PP). \square

The next lemma starts our analysis of gradient descent:

Lemma 2.8. *Let $H = \nabla G$ for $G \in C(X)$ such that ∇G is L -Lipschitz. Take $M_{i+1} \equiv I$ and $V'_{i+1}(u) := \tau_i(\nabla G(u^i) - \nabla G(u))$ with $W_{i+1} = \tau_i I$ as well as $Z_{i+1} \equiv \phi_i I$ for some $\tau_i, \phi_i > 0$. Then (CI \sim) holds with $U_{i+1} = U$ if*

$$(i) \quad \phi_i = \phi \text{ is constant, } \tau_i L < 2, \text{ and } \Delta_{i+1}(u^*; u) := -\phi_i(1 - \tau_i L/2) \|u - u^i\|^2/2.$$

If G is strongly convex with factor $\gamma > 0$, alternatively:

$$(ii) \quad \tau_0 L^2 < \gamma, \phi_{i+1} := \phi_i + \phi_i \tau_i (\gamma - \tau_i L^2), \tau_i := \phi_i^{-1/2}, \text{ and } \Delta_{i+1}(u^*; u) = 0.$$

Moreover, V_{i+1} satisfies (CL) under the above constraints on τ_i .

Proof. The satisfaction (CL) is immediate from the continuity of ∇G and the boundedness of τ_i . For the rest, we start by expanding the condition (CI*) as

$$(2.7) \quad \frac{\phi_i}{2} \|u - u^i\|^2 + \frac{\phi_i - \phi_{i+1}}{2} \|u - u^*\|^2 + \phi_i \tau_i \langle \nabla G(u^i) - \nabla G(u^*), u - u^* \rangle \geq -\Delta_{i+1}(u^*; u).$$

(i) Lipschitz gradient implies L^{-1} -co-coercivity ([1], see also Appendix B)

$$(2.8) \quad \langle \nabla G(u') - \nabla G(u), u' - u \rangle \geq L^{-1} \|\nabla G(u') - \nabla G(u)\|^2 \quad \text{for all } u, u'.$$

Now (2.7) follows after we use (2.8) and Cauchy's inequality to estimate

$$(2.9) \quad \begin{aligned} \langle \nabla G(u^i) - \nabla G(u^*), u - u^* \rangle &= \langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle \\ &\quad + \langle \nabla G(u^i) - \nabla G(u^*), u - u^i \rangle \geq -\frac{L}{4} \|u - u^i\|^2. \end{aligned}$$

(ii) We estimate

$$\begin{aligned} \langle \nabla G(u^i) - \nabla G(u^*), u - u^* \rangle &= \langle \nabla G(u) - \nabla G(u^*), u - u^* \rangle + \langle \nabla G(u^i) - \nabla G(u), u - u^* \rangle \\ &\geq \frac{\gamma}{2} \|u - u^*\|^2 - \frac{1}{2\tau_i} \|u - u^i\|^2 - \frac{\tau_i L^2}{2} \|u - u^*\|^2. \end{aligned}$$

Inserting this into (2.7), we see that (CI*) holds with $\Delta_{i+1}(u^*; u) = 0$ if

$$(2.10) \quad \phi_i + \phi_i \tau_i (\gamma - \tau_i L^2) \geq \phi_{i+1}.$$

Clearly our choice of $\{\tau_i\}_{i \in \mathbb{N}}$ is non-increasing. Therefore, (2.10) follows from the initialisation condition $\tau_0 L^2 < \gamma$ and the update rule $\phi_{i+1} := \phi_i + \phi_i \tau_i (\gamma - \tau_i L^2)$. \square

Example 2.3 (Gradient descent). Taking $\tau_i = \tau$ constant in Lemma 2.8, (PP) reads

$$0 = \tau \nabla G(u^i) + u^{i+1} - u^i.$$

This is the gradient descent method. Direct application of Lemma 2.8(i) with $u = u^{i+1}$ and $u^* = \hat{u}$ together with Corollary 2.2 and Proposition 2.5 now verifies the well-known weak convergence of the method when $\tau L < 2$.

Observe that $V_{i+1} = \nabla Q_{i+1}$ for

$$Q_{i+1}(u) := \frac{1}{2} \|u - u^i\|^2 + \tau [G(u^i) + \langle \nabla G(u^i), u - u^i \rangle - G(u)].$$

Each step of (PP) therefore minimises the *surrogate objective* [9]

$$(2.11) \quad u \mapsto G(u) + \tau^{-1} Q_{i+1}(u).$$

The function Q_{i+1} on one hand penalises long steps, and on the other hand allows longer steps when the local linearisation error is large. In this example, Q_{i+1} is, in fact, a Bregman distance. Proximal point methods based on general Bregman distances in place of the squared norm are studied in, e.g., [5, 8, 13, 14].

Example 2.4 (Acceleration of gradient descent). Continuing from [Example 2.3](#), if G is strongly convex, we may use the acceleration scheme in [Lemma 2.8\(ii\)](#). Similarly to [Example 2.1](#), ϕ_N is of the order $\Theta(N^2)$. Therefore, [Corollary 2.2](#) and [Proposition 2.4](#) show the convergence of $\|u^N - \widehat{u}\|^2$ to zero at the rate $O(1/N^2)$.

Example 2.5 (Forward–backward splitting). Let $H = \nabla G + \partial F$ for $G, F \in C(X)$ with ∇G Lipschitz. Taking M_{i+1} , W_{i+1} , and V'_{i+1} as in [Example 2.3](#), (PP) becomes

$$0 \in \tau_i \partial F(u^{i+1}) + \tau_i \nabla G(u^i) + u^{i+1} - u^i.$$

This is the forward–backward splitting method

$$u^{i+1} := (I + \tau_i \partial F)^{-1}(u^i - \tau_i \nabla G(u^i)).$$

By [Lemma 2.7](#), convergence and acceleration work exactly as for gradient descent in [Examples 2.3](#) and [2.4](#). If F is strongly convex with factor γ_F , we can introduce the additional term $\frac{\gamma_F}{2} \|u^{i+1} - \widehat{u}\|^2$ to (2.7). This will improve (2.10) to allow $\phi_{i+1} := \phi_i + \phi_i \tau_i (\gamma + \gamma_F - \tau_i L)$. Alternatively, it would be possible to choose ϕ_i and τ_i to yield FISTA-style acceleration [2].

Example 2.6 (Douglas–Rachford splitting). Let $A, B : U \rightrightarrows U$ be monotone operators. Consider the problem of finding \widehat{u} with $0 \in A(\widehat{u}) + B(\widehat{u})$. For $\lambda > 0$, let

$$(2.12) \quad \begin{aligned} H(u, v) &:= \begin{pmatrix} \lambda B(u) + u - v \\ \lambda A(u) + v - u \end{pmatrix}, & M_{i+1} &:= \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix}, \quad \text{and} \\ \widetilde{H}_{i+1}(u, v) &:= \begin{pmatrix} \lambda B(u^{i+1}) + u^{i+1} - v^i \\ \lambda A(u^{i+1} + v^{i+1} - v^i) + v^i - u^{i+1} \end{pmatrix}. \end{aligned}$$

Then $0 \in A(\widehat{u}) + B(\widehat{u})$ if and only if $0 \in H(\widehat{u}, \widehat{v})$, where $\widehat{v} \in (\widehat{u} - \lambda A(\widehat{u})) \cap (\widehat{u} + \lambda B(\widehat{u}))$. The algorithm (PP \sim) becomes the Douglas–Rachford splitting [10]

$$\begin{aligned} u^{i+1} &:= (I + \lambda B)^{-1}(v^i), \\ v^{i+1} &:= v^i + (I + \lambda A)^{-1}(2u^{i+1} - v^i) - u^{i+1}. \end{aligned}$$

We work with (PP \sim) since in (PP), V'_{i+1} would have to be set-valued. If A and B are maximal monotone, the variables $\{v^i\}_{i \in \mathbb{N}}$ converge weakly to \widehat{v} . Again, it is possible to devise

acceleration schemes under strong monotonicity [?, see, e.g.,] bredies2016accelerated.

Proof of convergence. Write $\bar{u}^i := (u^i, v^i)$ and $\widehat{u} := (\widehat{u}, \widehat{v})$. Observe that

$$u^{i+1} - v^{i+1} =: q^{i+1} \in \lambda A(u^{i+1} - v^{i+1} - v^i) \quad \text{and} \quad \widehat{u} - \widehat{v} =: \widehat{q} \in \lambda A(\widehat{u}).$$

Using the monotonicity of A and B , with $Z_{i+1} := I$, we have

$$\begin{aligned} \langle \widetilde{H}_{i+1}(\bar{u}^{i+1}), Z_{i+1}^*(\bar{u}^{i+1} - \widehat{u}) \rangle &\subset \langle \widetilde{H}_{i+1}(\bar{u}^{i+1}) - H(\widehat{u}), Z_{i+1}^*(\bar{u}^{i+1} - \widehat{u}) \rangle \\ &= \lambda \langle B(u^{i+1}) - B(\widehat{u}), u^{i+1} - \widehat{u} \rangle + \lambda \langle q^{i+1} - \widehat{q}, v^{i+1} - \widehat{v} \rangle \\ &\quad + \langle u^{i+1} - v^i, (u^{i+1} - v^{i+1}) - (\widehat{u} - \widehat{v}) \rangle \\ &= \lambda \langle B(u^{i+1}) - B(\widehat{u}), u^{i+1} - \widehat{u} \rangle + \lambda \langle q^{i+1} - \widehat{q}, u^{i+1} + v^{i+1} - v^i - \widehat{v} \rangle \geq 0. \end{aligned}$$

Thus (CI⁻) holds with $\Delta_{i+1}(\widehat{u}) := -\frac{1}{2}\|\bar{u}^{i+1} - \bar{u}^i\|_{Z_{i+1}M_{i+1}}^2$. Using (2.12) and the weak-to-strong outer semicontinuity of A and B (see Lemma A.1), we easily verify (CL). Weak convergence now follows from Theorem 2.1 and Proposition 2.5. \square

2.4 EXAMPLES OF SECOND-ORDER METHODS

Lemma 2.9. *Let $H = \nabla G$ for $G \in C^2(U)$. Take*

$$V_{i+1}(u) := \nabla^2 G(u^i)(u - u^i) + \nabla G(u^i) - \nabla G(u), \quad \text{and} \quad W_{i+1} := I$$

If $\nabla^2 G(u^) > 0$, then (CI^{*}) holds for u^i close enough to u^* with $\Delta_{i+1}(u, u^*) = 0$ and $Z_N M_N = \kappa^N \nabla^2 G(u^*)$ for some $\kappa > 1$.*

Proof. We set $M_{i+1} := \nabla^2 G(u^*)$ and $Z_{i+1} := \phi_i I$ for some $\phi_i > 0$. Then $G \in C^2(X)$ implies that $Z_{i+1} M_{i+1} = \phi_i \nabla^2 G(u^*)$ is self-adjoint. The condition (CI^{*}) reads

$$(2.13) \quad \frac{1}{2}\|u - u^i\|_{\phi_i \nabla^2 G(u^*)}^2 + \frac{1}{2}\|u - u^*\|_{(\phi_i - \phi_{i+1}) \nabla^2 G(u^*)}^2 + \phi_i D_{i+1} \geq -\Delta_{i+1}(u^*; u),$$

where

$$D_{i+1} := \langle \nabla G(u^i) - \nabla G(u^*) + (\nabla^2 G(u^i) - \nabla^2 G(u^*))(u - u^i), u - u^* \rangle.$$

By the fundamental theorem of calculus, there exists ζ^i between u^i and u^* with

$$D_{i+1} = \langle \nabla^2 G(\zeta^i)(u^i - u^*), u - u^* \rangle + \langle (\nabla^2 G(u^i) - \nabla^2 G(u^*))(u - u^i), u - u^* \rangle.$$

Using the three-point formula (2.5) and Cauchy's inequality we therefore obtain

$$\begin{aligned} D_{i+1} &= \frac{1}{2}\|u - u^*\|_{\nabla^2 G(u^i) - \nabla^2 G(u^*)}^2 - \frac{1}{2}\|u - u^i\|_{\nabla^2 G(u^*)}^2 + \frac{1}{2}\|u^i - u^*\|_{\nabla^2 G(u^*)}^2 \\ &\quad + \langle [\nabla^2 G(\zeta^i) - \nabla^2 G(u^i)](u^i - u^*), u - u^* \rangle \\ &\geq \frac{1}{2}\|u - u^*\|_{\nabla^2 G(u^i) - \nabla^2 G(u^*) - A_i}^2 - \frac{1}{2}\|u - u^i\|_{\nabla^2 G(u^*)}^2 \end{aligned}$$

for

$$A_i := [\nabla^2 G(\zeta^i) - \nabla^2 G(u^i)][\nabla^2 G(u^*)]^{-1}[\nabla^2 G(\zeta^i) - \nabla^2 G(u^i)].$$

Inserting this estimate into (2.13), we deduce that we can take $\Delta_{i+1}(u^*; u) = 0$ if

$$2\phi_i \nabla^2 G(u^i) - \phi_{i+1} \nabla^2 G(u^*) \geq \phi_i A_i.$$

Since $G \in C^2(U)$, and $\nabla^2 G(u^*) > 0$, locally near u^* , we can ensure $A_i \leq \epsilon \nabla^2 G(\zeta^{i+1})$ and $\nabla^2 G(u^i) \geq [\kappa/2 + \epsilon/2] \nabla^2 G(u^*)$ for some $\kappa > 1$ and $\epsilon > 0$. Thus it remains to satisfy

$$(1 + \epsilon)\kappa\phi_i - \phi_{i+1} \geq \phi_i \epsilon \kappa.$$

This holds when $\phi_{i+1} = \kappa\phi_i$. Taking $\phi_0 = 1$, thus $Z_N M_N \geq \kappa^N \nabla^2 G(u^*)$. \square

Example 2.7 (Newton's method). Suppose $H = \nabla G$ for $G \in C^2(U)$. Take V_{i+1} and W_{i+1} as in Lemma 2.9. Then (PP) reads

$$0 = \nabla G(u^i) + \nabla^2 G(u^i)(u^{i+1} - u^i).$$

This is Newton's method. By Lemma 2.9, Corollary 2.2, and Proposition 2.4, we obtain linear convergence if $\nabla^2 G(\bar{u}) > 0$.

Observe that now $V_{i+1}(u)$ is the gradient of

$$Q_{i+1}(u) := G(u^i) + \langle \nabla G(u^i), u - u^i \rangle + \frac{1}{2} \|u - u^i\|_{\nabla^2 G(u^i)}^2 - G(u).$$

In the surrogate objective (2.11), this allows longer steps when the second-order Taylor expansion under-approximates, and forces shorter steps when it over-approximates.

Example 2.8 (Proximal Newton's method). Similarly to Example 2.5, let $H = \nabla G + \partial F$ for $G \in C^2(X)$, and $F \in C(X)$. Taking M_{i+1} , W_{i+1} , and V'_{i+1} as in Lemma 2.9, (PP) becomes

$$0 \in \partial F(u^{i+1}) + \nabla G(u^i) + \nabla^2 G(u^i)(u^{i+1} - u^i).$$

This is the proximal Newton's method [?, see, e.g.,]lee2014proximal

$$u^{i+1} := (I + [\nabla^2 G(u^i)]^{-1} \partial F)^{-1}(u^i - [\nabla^2 G(u^i)]^{-1} \nabla G(u^i)),$$

where $(I + A^{-1} \partial F)^{-1}(v)$ solves $\min_u \frac{1}{2} \|u - v\|_A^2 + F(u)$. By Lemma 2.7, convergence and acceleration work exactly as for Newton's method in Example 2.7.

2.5 CONNECTIONS TO FIXED POINT THEOREMS

We demonstrate connections of our approach to established fixed point theorems.

Example 2.9 (Browder's fixed point theorem [4]). Let $T : U \rightarrow U$ be α -averaged, that is $T = (1 - \alpha)J + \alpha I$ for some non-expansive J and $\alpha \in (0, 1)$. Suppose there exists a fixed point

$\widehat{u} = T(\widehat{u})$. Then $u^i \rightarrow u^*$ for some fixed point u^* of T .

Proof of Browder's fixed point theorem. Let us set $H(u) := T(u) - u$, as well as $Z_{i+1} := W_{i+1} := M_{i+1} := I$ and $V'_{i+1}(u) := T(u^i) + u^i - T(u) - u$. We have

$$(2.14) \quad \widetilde{H}_{i+1}(u^{i+1}) := W_{i+1}H(u^{i+1}) + V'_{i+1}(u^{i+1}) = T(u^i) + u^i - 2u^{i+1} = u^i - u^{i+1},$$

where the last step follows by observing from the previous steps that (PP) says $u^{i+1} = T(u^i)$. The expression (2.14) easily gives (CL), and reduces (CI $\widetilde{}$) to

$$\frac{1}{2}\|u^{i+1} - u^i\|^2 + \langle u^i - u^{i+1}, u^{i+1} - \widehat{u} \rangle \geq -\Delta_{i+1}(\widehat{u}).$$

Using $u^{i+1} = T(u^i)$ and $\widehat{u} = T(\widehat{u})$, and taking $\beta > 0$, (CI $\widetilde{}$) therefore holds for

$$(2.15) \quad \Delta_{i+1}(\widehat{u}) = \frac{\alpha + 2\beta - 1}{2(1 - \alpha)}\|u^{i+1} - u^i\|^2$$

provided

$$0 \leq D := \frac{\beta}{1 - \alpha}\|T(u^i) - u^i\|^2 + \langle u^i - \widehat{u} - (T(u^i) - T(\widehat{u})), T(u^i) - T(\widehat{u}) \rangle.$$

Using the α -averaged property and $\widehat{u} = J(\widehat{u})$, we expand

$$\begin{aligned} \frac{D}{1 - \alpha} &= \beta\|J(u^i) - u^i\|^2 + \langle u^i - \widehat{u} - J(u^i) + J(\widehat{u}), (1 - \alpha)(J(u^i) - J(\widehat{u})) + \alpha(u^i - \widehat{u}) \rangle \\ &= (\alpha + \beta)\|u^i - \widehat{u}\|^2 + (\beta + \alpha - 1)\|J(u^i) - J(\widehat{u})\|^2 - (2\alpha + 2\beta - 1)\langle J(u^i) - J(\widehat{u}), u^i - \widehat{u} \rangle. \end{aligned}$$

We take $\beta := \max\{0, 1/2 - \alpha\}$. Then $2\alpha + 2\beta \geq 1$. Cauchy's inequality and non-expansivity of J thus give

$$\frac{D}{1 - \alpha} \geq \frac{1}{2}\|u^i - \widehat{u}\|^2 - \frac{1}{2}\|J(u^i) - J(\widehat{u})\|^2 \geq 0.$$

This verifies (CI $\widetilde{}$). From (2.15), $\Delta_{i+1}(\widehat{u}) \leq -\frac{1}{2}\min\{1, \alpha/(1 - \alpha)\}\|u^{i+1} - u^i\|^2$. We now obtain the claimed convergence from Corollary 2.2 and Proposition 2.5. \square

Remark 2.10. The preconditioner $V_{i+1}(u) = T(u^i) - T(u)$ is a T -based "distance", which is not obviously a Bregman distance.

3 SADDLE POINT PROBLEMS

With $K \in \mathcal{L}(X; Y)$, $G \in C(X)$ and $F^* \in C(Y)$ on Hilbert spaces X and Y , we now wish to solve (S). The first-order necessary optimality conditions can be written

$$(OC) \quad -K^*\widehat{y} \in \partial G(\widehat{x}), \quad \text{and} \quad K\widehat{x} \in \partial F^*(\widehat{y}).$$

Setting $U := X \times Y$ and introducing the variable splitting notation $u = (x, y)$, $\widehat{u} = (\widehat{x}, \widehat{y})$, etc., this succinctly be written as $0 \in H(\widehat{u})$ in terms of the operator

$$(3.1) \quad H(u) := \begin{pmatrix} \partial G(x) + K^*y \\ \partial F^*(y) - Kx \end{pmatrix}.$$

In this section, concentrating on this specific H , we specialise the theory of [Section 2.2](#) to saddle point problems. Throughout, for some primal and dual step length and testing operators $T_i, \Phi_i \in \mathcal{L}(X; X)$, and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{L}(Y; Y)$, we take

$$(3.2) \quad W_{i+1} := \begin{pmatrix} T_i & 0 \\ 0 & \Sigma_{i+1} \end{pmatrix}, \quad \text{and} \quad Z_{i+1} := \begin{pmatrix} \Phi_i & 0 \\ 0 & \Psi_{i+1} \end{pmatrix}.$$

To work with arbitrary step length operators, which will be necessary for stochastic algorithms in [Section 4.5](#), as well as the partially accelerated algorithms of [25], we will need abstract forms of partial strong monotonicity of G and F^* . As a first step, we take subsets of operators

$$\mathcal{T} \subset \mathcal{L}(X; X), \quad \text{and} \quad \mathcal{S} \subset \mathcal{L}(Y; Y).$$

We suppose that ∂G is *partially (strongly) \mathcal{T} -monotone*, which we take to mean

$$(G\text{-PM}) \quad \langle \partial G(x') - \partial G(x), x' - x \rangle_{\widetilde{T}} \geq \|x' - x\|_{\widetilde{T}}^2, \quad (x, x' \in X; \widetilde{T} \in \mathcal{T})$$

for some linear operator $0 \leq \Gamma \in \mathcal{L}(X; X)$. The operator $\widetilde{T} \in \mathcal{T}$ acts as a testing operator. Similarly, we assume that ∂F^* is *\mathcal{S} -monotone* in the sense

$$(F^*\text{-PM}) \quad \langle \partial F^*(y') - \partial F^*(y), y' - y \rangle_{\widetilde{S}} \geq 0 \quad (y, y' \in Y; \widetilde{S} \in \mathcal{S}).$$

Assuming G to satisfy (G-PM) for Γ and F^* to satisfy (F*-PM), we also introduce

$$\Xi_{i+1}(\Gamma) := \begin{pmatrix} 2T_i\Gamma & 2T_iK^* \\ -2\Sigma_{i+1}K & 0 \end{pmatrix},$$

which is an operator measure of strong monotonicity of H .

Example 3.1 (Block-separable structure, monotonicity). Let P_1, \dots, P_m be projection operators in X with $\sum_{j=1}^m P_j = I$ and $P_j P_i = 0$ if $i \neq j$. Suppose $G_1, \dots, G_m \in C(X)$ are (strongly) convex with factors $\gamma_1, \dots, \gamma_m \geq 0$. Then (G-PM) holds with $\Gamma = \sum_{j=1}^m \gamma_j P_j$ for

$$(3.3) \quad G(x) = \sum_{j=1}^m G_j(P_j x), \quad \text{and} \quad \mathcal{T} = \left\{ T := \sum_{j \in S} t_j P_j \mid t_j > 0, S \subset \{1, \dots, m\} \right\}.$$

3.1 ESTIMATES

Using the (strong) \mathcal{T} -monotonicity of ∂G , the next lemma simplifies [Corollary 2.2](#) for H given by (3.1). We introduce $\widetilde{\Gamma} = \Gamma$ to facilitate later gap estimates that will require the conditions in the lemma to hold for $\widetilde{\Gamma} = \Gamma/2$ instead of $\widetilde{\Gamma} = \Gamma$.

Theorem 3.1. *Let us be given $K \in \mathcal{L}(X; Y)$, $G \in C(X)$, and $F^* \in C(Y)$ on Hilbert spaces X and Y . Suppose G satisfies (G-PM) for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. For each $i \in \mathbb{N}$, let $T_i, \Phi_i \in \mathcal{L}(X; X)$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{L}(Y; Y)$ be such that $\Phi_i T_i \in \mathcal{T}$. Also take $V'_{i+1} : X \times Y \rightarrow X \times Y$, and $M_{i+1} \in \mathcal{L}(X \times Y; X \times Y)$. Let H given by (3.1), Z_{i+1} and W_{i+1} by (3.2), and V_{i+1} by (2.3). Suppose (PP) is solvable, and denote the iterates by $u^i = (x^i, y^i)$. Then (CI), (CI $\bar{\Gamma}$) and (DI) hold if $Z_{i+1}M_{i+1}$ is self-adjoint, and for $\tilde{\Gamma} = \Gamma$ we have*

$$\begin{aligned}
(\text{CI-}\Gamma) \quad & \underbrace{\frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2}_{\text{step length in local metric}} + \underbrace{\frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}(\Xi_{i+1}(\tilde{\Gamma})+M_{i+1})-Z_{i+2}M_{i+2}}^2}_{\text{linear preconditioner update discrepancy}} \\
& + \underbrace{\langle \partial F^*(y^{i+1}) - \partial F^*(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}}}_{\text{variably useful remainder from } H} + \underbrace{\langle V'_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}}}_{\text{from non-linear preconditioner}} \\
& \geq -\Delta_{i+1}(\widehat{u}).
\end{aligned}$$

Proof. First of all, we observe that (CI- Γ) implies

$$\begin{aligned}
(3.4) \quad & \frac{1}{2} \|u^{i+1} - u^i\| + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}(\Xi_{i+1}(0)+M_{i+1})-Z_{i+2}M_{i+2}} \\
& + \langle \partial G(x^{i+1}) - \partial G(\widehat{x}), x^{i+1} - \widehat{x} \rangle_{\Phi_i T_i} + \langle \partial F^*(y^{i+1}) - \partial F^*(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}} \\
& + \langle V'_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}} \geq -\Delta_{i+1}(\widehat{u}).
\end{aligned}$$

Here pay attention to the fact that (3.4) employs $\Xi_{i+1}(0)$ while (CI- Γ) employs $\Xi_{i+1}(\tilde{\Gamma})$. If we show that (CI) follows from (3.4), then (CI $\bar{\Gamma}$) and (DI) follow from Corollary 2.2. Indeed, using the expansion

$$Z_{i+1}W_{i+1} = \begin{pmatrix} \Phi_i T_i & 0 \\ 0 & \Psi_{i+1}\Sigma_{i+1} \end{pmatrix},$$

we expand for any $\tilde{u} = (\tilde{x}, \tilde{y})$ that

$$\begin{aligned}
& \langle Z_{i+1}W_{i+1}(H(u^{i+1}) - H(\tilde{u})), u^{i+1} - \tilde{u} \rangle \\
& = \langle \partial G(x^{i+1}) - \partial G(\tilde{x}), x^{i+1} - \tilde{x} \rangle_{\Phi_i T_i} + \langle \partial F^*(y^{i+1}) - \partial F^*(\tilde{y}), y^{i+1} - \tilde{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}} \\
& + \langle \Phi_i T_i K^*(y^{i+1} - \tilde{y}), x^{i+1} - \tilde{x} \rangle - \langle \Psi_{i+1}\Sigma_{i+1}K(x^{i+1} - \tilde{x}), y^{i+1} - \tilde{y} \rangle.
\end{aligned}$$

With the help of $\Xi_{i+1}(0)$ we then obtain

$$\begin{aligned}
& \langle H(u^{i+1}) - H(\tilde{u}), u^{i+1} - \tilde{u} \rangle_{Z_{i+1}W_{i+1}} \geq \frac{1}{2} \|u^{i+1} - \tilde{u}\|_{Z_{i+1}\Xi_{i+1}(0)} \\
& + \langle \partial G(x^{i+1}) - \partial G(\tilde{x}), x^{i+1} - \tilde{x} \rangle_{\Phi_i T_i} + \langle \partial F^*(y^{i+1}) - \partial F^*(\tilde{y}), y^{i+1} - \tilde{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}}.
\end{aligned}$$

Inserting this into (3.4), we obtain (CI). \square

3.2 EXAMPLES OF PRIMAL-DUAL METHODS

We now look at several known methods for the saddle point problem (S).

Example 3.2 (The primal–dual method of Chambolle and Pock [6]). This method consists of iterating the system

$$(3.5a) \quad x^{i+1} := (I + \tau_i \partial G)^{-1}(x^i - \tau_i K^* y^i),$$

$$(3.5b) \quad \bar{x}^{i+1} := \omega_i(x^{i+1} - x^i) + x^{i+1},$$

$$(3.5c) \quad y^{i+1} := (I + \sigma_{i+1} \partial F^*)^{-1}(y^i + \sigma_{i+1} K \bar{x}^{i+1}).$$

In the basic version of the algorithm, $\omega_i = 1$, $\tau_i \equiv \tau_0 > 0$, and $\sigma_i \equiv \sigma_0 > 0$, assuming the step length parameters to satisfy

$$(3.6) \quad \tau_0 \sigma_0 \|K\|^2 < 1.$$

The iterates convergence weakly, and the method has $O(1/N)$ rate for the ergodic duality gap, to which we will return in [Section 4](#). If G is strongly convex with factor γ , we may take $\tilde{\gamma} \in (0, \gamma]$, and accelerate

$$(3.7) \quad \omega_i := 1/\sqrt{1 + 2\tilde{\gamma}\tau_i}, \quad \tau_{i+1} := \tau_i \omega_i, \quad \text{and} \quad \sigma_{i+1} := \sigma_i / \omega_i.$$

This yields $O(1/N^2)$ convergence of $\|x^N - \hat{x}\|^2$ to zero.

Proof of convergence of iterates. We formulate the method in our proximal point framework following [25, 12] by taking as the preconditioner

$$M_{i+1} = \begin{pmatrix} I & -\tau_i K^* \\ -\sigma_i K & I \end{pmatrix} \quad \text{and} \quad V'_{i+1} = 0.$$

As the step length and testing operators we take $T_i = \tau_i I$, $\Sigma_{i+1} = \sigma_{i+1} I$, $\Phi_i = \phi_i I$, $\Psi_{i+1} = \psi_{i+1} I$. We also write $\tilde{\Gamma} := \tilde{\gamma} I$. Taking $\Delta_{i+1}(\hat{u}) := -\frac{1}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$, we reduce (CI- Γ) to

$$(3.8) \quad \frac{1}{2}\|u^{i+1} - \hat{u}\|_{D_{i+2}}^2 \geq 0 \quad \text{for} \quad D_{i+2} := Z_{i+1}(\Xi_{i+1}(\tilde{\Gamma}) + M_{i+1}) - Z_{i+2}M_{i+2}.$$

We may expand

$$(3.9a) \quad Z_{i+1}M_{i+1} = \begin{pmatrix} \phi_i I & -\phi_i \tau_i K^* \\ -\psi_{i+1} \sigma_i K & \psi_{i+1} I \end{pmatrix}, \quad \text{and}$$

$$(3.9b) \quad D_{i+2} = \begin{pmatrix} (\phi_i(1 + 2\tilde{\gamma}\tau_i) - \phi_{i+1})I & (\phi_i \tau_i + \phi_{i+1} \tau_{i+1})K^* \\ (\psi_{i+2} \sigma_{i+1} - 2\psi_{i+1} \sigma_{i+1} - \psi_{i+1} \sigma_i)K & (\psi_{i+1} - \psi_{i+2})I \end{pmatrix}.$$

We have $\|\cdot\|_{D_{i+2}} = 0$ (but not $D_{i+2} = 0$, as the former depends on the off-diagonals cancelling out), and $Z_{i+1}M_{i+1}$ is self-adjoint, if for some constant ψ we take

$$(3.10) \quad \phi_{i+1} := \phi_i(1 + 2\tilde{\gamma}\tau_i), \quad \tau_i := \phi_i^{-1/2}, \quad \sigma_i := \phi_i \tau_i / \psi, \quad \text{and} \quad \psi_{i+1} := \psi.$$

This gives the acceleration scheme (3.7). Moreover, for any $\delta \in (0, 1)$ holds

$$(3.11) \quad Z_{i+1}M_{i+1} \geq \begin{pmatrix} \delta \phi_i I & 0 \\ 0 & \psi I - (1 - \delta)^{-1} K K^* \end{pmatrix}.$$

Thus $Z_{i+1}M_{i+1} \geq 0$ if $\psi \geq (1 - \delta)^{-1}\|K\|^2$. By (3.10), $\sigma_i\tau_i = 1/\psi$. Since this fixes the ratio of σ_i to τ_i , we need to take $\psi := 1/(\sigma_0\tau_0)$ as well as $\delta := 1 - \sigma_0\tau_0\|K\|^2$. Through the positivity of δ , we recover the initialisation condition (3.6).

Theorem 3.1 and **Proposition 2.5** show weak convergence of the iterates without a rate. If G is strongly convex with factor $\gamma \geq 0$, so that also $\tilde{\gamma} > 0$, the results in [6, 25] show that τ_N is of the order $O(1/N)$, and consequently ϕ_N is of the order $\Theta(N^2)$. By **Proposition 2.4**, $\|x^N - \hat{x}\|^2$ converges to zero at the rate $O(1/N^2)$. \square

Example 3.3 (Alternating Directions Method of Multipliers, briefly). The classical ADMM [11] and Douglas–Rachford splitting [10] are known to be related to the Chambolle–Pock method; in fact the Chambolle–Pock method is a preconditioned ADMM [6]. From [3, Section 5], we can deduce that compared to the Chambolle–Pock method, the ADMM merely has the sign of K reversed in

$$M_{i+1} = \begin{pmatrix} I & \tau_i K \\ \sigma_i K & I \end{pmatrix}.$$

Taking $\tau_i = \tau_0$ and $\sigma_i = \sigma_0$ constant and satisfying (3.6), the iterates converge weakly. Acceleration can provide $O(1/N)$ convergence of $\|x^N - \hat{x}\|^2$.

Proof of convergence. Following **Example 3.2**, we now expand

$$D_{i+2} = \begin{pmatrix} (\phi_i(1 + 2\tilde{\gamma}\tau_i) - \phi_{i+1})I & (3\phi_i\tau_i - \phi_{i+1}\tau_{i+1})K^* \\ (\psi_{i+1}\sigma_i - 2\psi_{i+1}\sigma_{i+1} - \psi_{i+2}\sigma_{i+1})K & (\psi_{i+1} - \psi_{i+2})I \end{pmatrix}.$$

This time $\|\cdot\|_{D_{i+2}} = 0$ and $Z_{i+1}M_{i+1}$ is self-adjoint if we take

$$(3.12) \quad \phi_{i+1} := \phi_i(1 + 2\tilde{\gamma}\tau_i), \quad \tau_{i+1} := \tau_i\phi_i/\phi_{i+1}, \quad \sigma_i := \phi_i\tau_i/\psi, \quad \text{and} \quad \psi_{i+1} := \psi.$$

If $\tilde{\gamma} = 0$, which corresponds to the standard ADMM with fixed step lengths, it is easy to retrace the steps of **Example 3.2** to prove weak convergence (without a rate). If $\tilde{\gamma} \neq 0$, we obtain $\phi_{N+1} = \phi_N + 2\tilde{\gamma}\tau_{N-1}\phi_{N-1} = \phi_N + 2\tilde{\gamma}\tau_0\phi_0 = \phi_0 + 2N\tilde{\gamma}\tau_0\phi_0$. Therefore, the acceleration scheme (3.12) only gives the rate $O(1/N)$. \square

Example 3.4 (Chambolle–Pock with a forward step). Suppose $G = G_0 + J$ with G (strongly) convex with factor $\gamma \geq 0$, and ∇J Lipschitz with factor L . (J does not have to be convex.) In [7], the Chambolle–Pock method was extended to take forward steps with respect to J . With everything else as in **Example 3.2**, take $V'_{i+1}(u) := (\tau_i(\nabla J(x^i) - \nabla J(x)), 0)$. Then (PP) can be rearranged as

$$(3.13) \quad x^{i+1} := (I + \tau_i\partial G_0)^{-1}(x^i - \tau_i\nabla J(x^i) - \tau_i K^* y^i),$$

$$(3.14) \quad \bar{x}^{i+1} := \omega_i(x^{i+1} - x^i) + x^{i+1},$$

$$(3.15) \quad y^{i+1} := (I + \sigma_{i+1}\partial F^*)^{-1}(y^i + \sigma_{i+1}K\bar{x}^{i+1}).$$

The method inherits the convergences properties of **Example 3.2** if we use the step length

update rules (3.7), and initialise $\tau_0, \sigma_0 > 0$ subject to (3.6), and

$$(3.16) \quad 0 < \theta := 1 - L\tau_0/(1 - \tau_0\sigma_0\|K\|^2).$$

Proof of convergence. With D_{i+2} as in (3.8), the condition (CI- Γ) becomes

$$(3.17) \quad \frac{1}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2}\|u^{i+1} - \widehat{u}\|_{D_{i+2}}^2 + \tau_i\phi_i\langle \nabla J(x^i) - \nabla J(\widehat{x}), x^{i+1} - \widehat{x} \rangle \geq -\Delta_{i+1}(\widehat{u}).$$

The rules (3.10) force $\|\cdot\|_{D_{i+2}} = 0$. Applying the estimate (2.9) to J , (3.17) becomes

$$\frac{1}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 - \frac{\tau_i\phi_i L}{4}\|x^{i+1} - x^i\|^2 \geq -\Delta_{i+1}(\widehat{u}).$$

We take $\Delta_{i+1}(\widehat{u}) = -\frac{\theta}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2$ for some $\theta > 0$, and deduce using Cauchy's inequality that this condition holds if

$$(1 - \theta)Z_{i+1}M_{i+1} \geq \tau_i\phi_i L \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix}.$$

Recalling (3.11), this is true if $(1 - \theta)\delta\phi_i \geq \tau_i\phi_i L$ and $\psi \geq (1 - \delta)^{-1}\phi_i\tau_i^2\|K\|^2$. Further recalling (3.10), and observing that $\{\tau_i\}$ is non-increasing, we only have to satisfy $(1 - \theta)(1 - \tau_0\sigma_0\|K\|^2) \geq L\tau_0$. Otherwise put, we obtain (3.16). \square

Example 3.5 (GIST). Suppose $G(x) = \frac{1}{2}\|f - Ax\|^2$, $\|A\| < \sqrt{2}$, and $\|K\| \leq 1$. Take

$$V'_{i+1}(u) := \begin{pmatrix} \nabla G(x^i) - \nabla G(x) \\ 0 \end{pmatrix}, \quad \text{and} \quad M_{i+1} := \begin{pmatrix} I & 0 \\ 0 & I - KK^* \end{pmatrix}.$$

With $T_i := I$ and $\Sigma_{i+1} := I$, we then obtain the Generalised Iterative Soft Thresholding (GIST) algorithm of [16]

$$\begin{aligned} y^{i+1} &:= (I + \partial F^*)^{-1}((I - KK^*)y^i + K(x^i - \nabla G(x^i))), \\ x^{i+1} &:= x^i - \nabla G(x^i) - K^*y^{i+1}. \end{aligned}$$

The iterates $\{x^i\}_{i \in \mathbb{N}}$ converge weakly to \widehat{x} .

Proof of convergence. Clearly $Z_{i+1}M_{i+1}$ is positive semi-definite self-adjoint. Also G satisfies (G-PM) with $\Gamma = A^*A$. If we take $\Phi_i = I$ and $\Psi_{i+1} = I$, then

$$D_{i+2} := Z_{i+1}(\Xi_{i+1}(\widetilde{\Gamma}) + M_{i+1}) - Z_{i+2}M_{i+2} = \begin{pmatrix} 2A^*A & 2K^* \\ -2K & 0 \end{pmatrix}.$$

Thus $\frac{1}{2}\|u\|_{D_{i+2}}^2 = \|x\|_{A^*A}^2$. Eliminating ∂F^* by monotonicity, (CI- Γ) thus holds if

$$\frac{1}{2}\|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \|x^{i+1} - \widehat{x}\|_{A^*A}^2 + \langle Z_{i+1}V'_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle \geq -\Delta_{i+1}(\widehat{u}).$$

Expanding and using $\|K\| < 1$, we see this to hold when

$$\frac{1}{2}\|x^{i+1} - x^i\|^2 + \|x^{i+1} - \widehat{x}\|_{A^*A}^2 + \langle A^*A(x^i - x^{i+1}), x^{i+1} - \widehat{x} \rangle \geq -\Delta_{i+1}(\widehat{u}).$$

Our assumption $\|A\| < \sqrt{2}$ guarantees $\frac{1}{2}(A^*A)^2 < A^*A$. Cauchy's inequality therefore shows that we can take $\Delta_{i+1} = -\frac{c}{2}\|x^{i+1} - x^i\|^2$ for some $c > 0$. Using [Theorem 3.1](#) and [Proposition 2.4](#), we obtain weak convergence. \square

4 THE ERGODIC DUALITY GAP AND STOCHASTIC METHODS

We now study the extension of the testing approach of [Section 2.2](#) to produce the convergence of an ergodic duality gap. Throughout this section, we are in the saddle point setup of [Section 3](#). In particular, H is as in [\(3.1\)](#), and the step length and testing operators W_{i+1} and Z_{i+1} as in [\(3.2\)](#).

4.1 PRELIMINARY GAP ESTIMATES

Our first lemma demonstrates how to obtain a ‘‘preliminary’’ gap $\mathcal{G}'_{i+1}(u)$ from H . If the step lengths and tests are scalar, $T_i = \tau_i I$, and $\Phi_i = \phi_i I$, etc., and satisfy $\tau_i \phi_i = \sigma_i \psi_{i+1}$, it is easy to bound this preliminary gap from below by $\tau_i \phi_i$ times the conventional duality gap

$$(4.1) \quad \mathcal{G}(x, y) := (G(x) + \langle \widehat{y}, Kx \rangle - F(\widehat{y})) - (G(\widehat{x}) + \langle y, K\widehat{x} \rangle - F^*(y)).$$

To do the same for more general step length operators, we will in [Section 4.2](#) introduce abstract notions of convexity that incorporate ergodicity and stochasticity.

Lemma 4.1. *Let us be given $K \in \mathcal{L}(X; Y)$, $G \in C(X)$, and $F^* \in C(Y)$ on Hilbert spaces X and Y . For each $i \in \mathbb{N}$, let $T_i, \Phi_i \in \mathcal{L}(X; X)$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{L}(Y; Y)$. Then for any $\widetilde{\Gamma} \in \mathcal{L}(X; X)$,*

$$(4.2) \quad \langle H(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} = \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) + \frac{1}{2}\|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Xi_{i+1}(\widetilde{\Gamma})}^2,$$

where the ‘‘preliminary gap’’

$$\begin{aligned} \mathcal{G}'_{i+1}(u; \widetilde{\Gamma}) := & \langle \partial G(x), x - \widehat{x} \rangle_{\Phi_i T_i} - \|x - \widehat{x}\|_{\Phi_i T_i \widetilde{\Gamma}}^2 + \langle \partial F^*(y), y - \widehat{y} \rangle_{\Psi_{i+1} \Sigma_{i+1}} \\ & - \langle \widehat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1} \Sigma_{i+1} K) \widehat{x} \rangle - \langle y, \Psi_{i+1} \Sigma_{i+1} K \widehat{x} \rangle + \langle \widehat{y}, KT_i^* \Phi_i^* x \rangle. \end{aligned}$$

Proof. Similarly to the proof of [Theorem 3.1](#), we have

$$\begin{aligned} \langle H(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} = & \langle \partial G(x^{i+1}), x^{i+1} - \widehat{x} \rangle_{\Phi_i T_i} + \langle \Phi_i T_i K^* y^{i+1}, x^{i+1} - \widehat{x} \rangle \\ & + \langle \partial F^*(y^{i+1}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1} \Sigma_{i+1}} - \langle \Psi_{i+1} \Sigma_{i+1} K x^{i+1}, y^{i+1} - \widehat{y} \rangle. \end{aligned}$$

A little bit of reorganisation gives (4.2). Indeed

$$\begin{aligned}
\langle H(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} &= \langle \partial G(x^{i+1}), x^{i+1} - \widehat{x} \rangle_{\Phi_i T_i} - \|x^{i+1} - \widehat{x}\|_{\Phi_i T_i \widetilde{\Gamma}}^2 \\
&\quad + \langle \partial F^*(y^{i+1}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}} + \|x^{i+1} - \widehat{x}\|_{\Phi_i T_i \widetilde{\Gamma}}^2 \\
&\quad + \langle y^{i+1} - \widehat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1}\Sigma_{i+1}K)(x^{i+1} - \widehat{x}) \rangle \\
&\quad - \langle \widehat{y}, (KT_i^* \Phi_i^* - \Psi_{i+1}\Sigma_{i+1}K)\widehat{x} \rangle \\
&\quad - \langle y^{i+1}, \Psi_{i+1}\Sigma_{i+1}K\widehat{x} \rangle + \langle \widehat{y}, KT_i^* \Phi_i^* x^{i+1} \rangle \\
&= \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Xi_{i+1}(\widetilde{\Gamma})}^2. \quad \square
\end{aligned}$$

The next lemma extends [Theorem 3.1](#) to estimate the preliminary gap.

Lemma 4.2. *Let us be given $K \in \mathcal{L}(X; Y)$, $G \in C(X)$, and $F^* \in C(Y)$ on Hilbert spaces X and Y . For each $i \in \mathbb{N}$, let $T_i, \Phi_i \in \mathcal{R}(\mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(\mathcal{L}(Y; Y))$, as well as $V'_{i+1} \in \mathcal{R}(X \times Y \rightarrow X \times Y)$ and $M_{i+1} \in \mathcal{R}(\mathcal{L}(X \times Y; X \times Y))$. Let H given by (3.1), Z_{i+1} and W_{i+1} by (3.2), and V_{i+1} by (2.3). Suppose (PP) is solvable, and denote the iterates by $u^i = (x^i, y^i)$. If $Z_{i+1}M_{i+1}$ is self-adjoint, and*

$$\begin{aligned}
(\text{CI-}\mathcal{G}) \quad \frac{1}{2} \|u^{i+1} - u^i\|_{Z_{i+1}M_{i+1}}^2 + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}(\Xi_{i+1}(\widetilde{\Gamma}) + M_{i+1}) - Z_{i+2}M_{i+2}}^2 \\
+ \langle V'_{i+1}(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}} \geq -\widetilde{\Delta}_{i+1}(\widehat{u})
\end{aligned}$$

for some $\widetilde{\Gamma} \in \mathcal{L}(X; X)$, then

$$(4.3) \quad \frac{1}{2} \|u^N - \widehat{u}\|_{Z_{N+1}M_{N+1}}^2 + \sum_{i=0}^{N-1} \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) \leq \frac{1}{2} \|u^0 - \widehat{u}\|_{Z_1M_1}^2 + \sum_{i=0}^{N-1} \widetilde{\Delta}_{i+1}(\widehat{u}) \quad (N \geq 1).$$

Proof. Inserting (4.2) into (CI- \mathcal{G}) proves (CI \sim) for $\Delta_{i+1}(\widehat{u}) := \widetilde{\Delta}_{i+1}(\widehat{u}) - \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma})$. Now we use [Theorem 2.1](#). \square

The problem with the above [Lemma 4.2](#) is that it loses ∂F from the condition (CI- \mathcal{G}) compared to (CI- Γ). Thus (CI- \mathcal{G}) can be more difficult to satisfy for particular preconditioners that are related to ∂F , such as the forward-backward splitting in [Example 2.5](#). Fortunately, there is a remedy: to study a one-sided gap that provides no indication of the convergence of the dual variable.

Lemma 4.3. *Let us be given $K \in \mathcal{L}(X; Y)$, $G \in C(X)$, and $F^* \in C(Y)$ on Hilbert spaces X and Y . For each $i \in \mathbb{N}$, let $T_i, \Phi_i \in \mathcal{R}(\mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(\mathcal{L}(Y; Y))$, as well as $V'_{i+1} \in \mathcal{R}(X \times Y \rightarrow X \times Y)$ and $M_{i+1} \in \mathcal{R}(\mathcal{L}(X \times Y; X \times Y))$. Let H given by (3.1), Z_{i+1} and W_{i+1} by (3.2), and V_{i+1} by (2.3). Suppose (PP) is solvable, and denote the iterates by $u^i = (x^i, y^i)$. If $Z_{i+1}M_{i+1}$ is self-adjoint, and (CI- Γ) holds for some $\widetilde{\Gamma} \in \mathcal{L}(X; X)$, then*

$$(4.4) \quad \frac{1}{2} \|u^N - \widehat{u}\|_{Z_{N+1}M_{N+1}}^2 + \sum_{i=0}^{N-1} \mathcal{G}'_{i+1}(x^{i+1}, \widehat{y}; \widetilde{\Gamma}) \leq \frac{1}{2} \|u^0 - \widehat{u}\|_{Z_1M_1}^2 + \sum_{i=0}^{N-1} \Delta_{i+1}(\widehat{u}) \quad (N \geq 1).$$

Proof. Let us write $(H_x(u), H_y(u)) := H(u)$. Then $0 \in H_y(\widehat{y})$. We may thus expand

$$\begin{aligned} \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) &= \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) - \langle H_y(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}} \\ &= \mathcal{G}'_{i+1}(u^{i+1}; \widetilde{\Gamma}) - \langle \partial F^*(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}} + \langle \Psi_{i+1}\Sigma_{i+1}K^*\widehat{x}, y^{i+1} - \widehat{y} \rangle \\ &= \mathcal{G}'_{i+1}(x^{i+1}, \widehat{y}; \widetilde{\Gamma}) + \langle \partial F^*(y^{i+1}) - \partial F^*(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}}. \end{aligned}$$

Inserting this into (4.2), we obtain

$$(4.5) \quad \begin{aligned} \langle H(u^{i+1}), u^{i+1} - \widehat{u} \rangle_{Z_{i+1}W_{i+1}} &= \mathcal{G}'_{i+1}(x^{i+1}, \widehat{y}; \widetilde{\Gamma}) + \frac{1}{2} \|u^{i+1} - \widehat{u}\|_{Z_{i+1}\Xi_{i+1}(\widetilde{\Gamma})}^2 \\ &\quad + \langle \partial F^*(y^{i+1}) - \partial F^*(\widehat{y}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1}\Sigma_{i+1}}. \end{aligned}$$

We write $\widetilde{\Delta}_{i+1}$ for the Δ_{i+1} for which (CI- Γ) holds. Inserting (4.5) into (CI- Γ) proves (CI- $\widetilde{\Gamma}$) for $\Delta_{i+1}(\widehat{u}) := \widetilde{\Delta}_{i+1}(\widehat{u}) - \mathcal{G}'_{i+1}(x^{i+1}, \widehat{y}; \widetilde{\Gamma})$. The rest follows from Theorem 2.1. \square

4.2 CONVERSION OF PRELIMINARY GAPS TO ERGODIC GAPS

The ‘‘preliminary gaps’’ are not as such very useful. To go further, the abstract monotonicity assumptions (G-PM) and (F*-PM) are not enough, and we need analogous convexity formulations. We formulate these conditions directly in the stochastic setting. Towards this end we introduce the following notation:

Definition 4.1. We write $x \in \mathcal{R}(X)$ if T is an \mathcal{T} -valued random variable: $x : \Omega \rightarrow X$ for some (in the present work fixed) probability space (Ω, \mathcal{O}) , where \mathcal{O} is a σ -algebra on Ω . We denote by \mathbb{E} the expectation with respect to a probability measure \mathbb{P} on Ω . As is common, we abuse notation and write $x = x(\omega)$ for the unknown random realisation $\omega \in \Omega$.

We refer to [23] for more details on measure-theoretic probability. From now on, we assume for all $N \geq 1$ that whenever $\widetilde{T}_i (:= \Phi_i T_i) \in \mathcal{R}(\mathcal{T})$ and $x^{i+1} \in \mathcal{R}(X)$ for each $i = 0, \dots, N-1$ with $\sum_{i=0}^{N-1} \mathbb{E}[\widetilde{T}_i] = I$, then for some $0 \leq \Gamma \in \mathcal{L}(X; X)$ holds

$$(G\text{-EC}) \quad G\left(\sum_{i=0}^{N-1} \mathbb{E}[\widetilde{T}_i^* x^{i+1}]\right) - G(\widehat{x}) \geq \sum_{i=0}^{N-1} \mathbb{E}\left[\langle \partial G(x^{i+1}), x^{i+1} - \widehat{x} \rangle_{\widetilde{T}_i} + \frac{1}{2} \|x^{i+1} - \widehat{x}\|_{\widetilde{T}_i \Gamma}^2\right].$$

Analogously, we assume for $\widetilde{\Sigma}_{i+1} (:= \Psi_{i+1}\Sigma_{i+1}) \in \mathcal{R}(\mathcal{S})$ and $y^{i+1} \in \mathcal{R}(Y)$ for each $i = 0, \dots, N-1$ with $\sum_{i=0}^{N-1} \mathbb{E}[\widetilde{\Sigma}_{i+1}] = I$ that

$$(F^*\text{-EC}) \quad F^*\left(\sum_{i=0}^{N-1} \mathbb{E}[\widetilde{\Sigma}_{i+1}^* y^{i+1}]\right) - F^*(\widehat{y}) \geq \sum_{i=0}^{N-1} \mathbb{E}\left[\langle \partial F^*(y^{i+1}), y^{i+1} - \widehat{y} \rangle_{\widetilde{\Sigma}_{i+1}}\right].$$

Example 4.1 (Block-separable structure, ergodic convexity). Let G and \mathcal{T} have the separable structure of Example 3.1. We claim that (G-EC) holds. Indeed, let us introduce $\widetilde{T}_i := \sum_{j=1}^m \widetilde{\tau}_{j,i} P_j \geq 0$, satisfying $\sum_{i=0}^{N-1} \mathbb{E}[\widetilde{\tau}_{j,i}] = 1$ for each $j = 1, \dots, m$. Splitting (G-EC) into separate inequalities over all $j = 1, \dots, m$, and using the strong convexity of G_j , we see (G-EC)

to be true if for all $j = 1, \dots, m$ holds

$$(4.6) \quad G_j \left(\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\tau}_{j,i} P_j x^{i+1}] \right) - G_j(P_j \widehat{x}) \geq \sum_{i=0}^{N-1} \mathbb{E} \left[\tilde{\tau}_i (G_j(P_j x^{i+1}) - G_j(P_j \widehat{x})) \right].$$

The right hand side can also be written as $\int_{\Omega^N} G_j(P_j x^i(\omega)) - G_j(P_j \widehat{x}) d\mu^N(i, \omega)$ for the measure $\mu^N := \tilde{\tau}_j \sum_{i=0}^{N-1} \delta_i \times \mathbb{P}$ on the domain $\Omega^N := \{0, \dots, N-1\} \times \Omega$. Using our assumption $\sum_{i=0}^{N-1} \mathbb{E}[\tilde{\tau}_{j,i}] = 1$, we deduce $\mu^N(\Omega^N) = 1$. An application of Jensen's inequality now shows (4.6). Therefore (G-EC) is satisfied.

We also assume that either

$$(CG) \quad \mathbb{E}[\Phi_i T_i] = \bar{\eta}_i I, \quad \text{and} \quad \mathbb{E}[\Psi_{i+1} \Sigma_{i+1}] = \bar{\eta}_i I, \quad (i \geq 1),$$

or

$$(CG_*) \quad \mathbb{E}[\Phi_i T_i] = \bar{\eta}_i I, \quad \text{and} \quad \mathbb{E}[\Psi_i \Sigma_i] = \bar{\eta}_i I, \quad (i \geq 1),$$

As will see in [Example 4.2](#), (CG_*) is satisfied by the accelerated Chambolle–Pock method of [Example 3.2](#). In our companion paper [24], we will however see that (CG) is required to develop doubly-stochastic methods.

With these, and the gap functional \mathcal{G} from (4.1), we derive the next two lemmas that are meant to be used in combination with either [Lemma 4.2](#) or [Lemma 4.3](#), to estimate the sum of the preliminary gaps therein. For this, the expectation needs to be taken in the estimates of the latter. All of these different combinations will be summarised in [Theorem 4.6](#) after the lemmas.

Lemma 4.4. *Suppose (G-EC), (F*-EC), and (CG) hold. Set*

$$(4.7) \quad \zeta_N := \sum_{i=0}^{N-1} \bar{\eta}_i,$$

and for $\{(x^i, y^i)\}_{i=1}^N \subset X \times Y$, define the ergodic sequences

$$(4.8) \quad \tilde{x}_N := \zeta_N^{-1} \mathbb{E} \left[\sum_{i=0}^{N-1} T_i^* \Phi_i^* x^{i+1} \right], \quad \text{and} \quad \tilde{y}_N := \zeta_N^{-1} \mathbb{E} \left[\sum_{i=0}^{N-1} \Sigma_{i+1}^* \Psi_{i+1}^* y^{i+1} \right].$$

Then

$$\sum_{i=0}^{N-1} \mathbb{E}[\mathcal{G}'_{i+1}(x^{i+1}, y^{i+1}; \Gamma/2)] \geq \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N).$$

Proof. Using (CG) , $(G-EC)$, and (F^*-EC) , we compute

$$\begin{aligned} \sum_{i=0}^{N-1} \mathbb{E}[\mathcal{G}'_{i+1}(x^{i+1}, y^{i+1}; \Gamma/2)] &= \sum_{i=0}^{N-1} \mathbb{E} \left[\langle \partial G(x^{i+1}), x^{i+1} - \widehat{x} \rangle_{\Phi_i T_i} \right. \\ &\quad \left. - \|x^{i+1} - \widehat{x}\|_{\Phi_i T_i \Gamma/2}^2 + \langle \partial F^*(y^{i+1}), y^{i+1} - \widehat{y} \rangle_{\Psi_{i+1} \Sigma_{i+1}} \right] \\ &\quad - \zeta_N \langle \tilde{y}_N, K \widehat{x} \rangle + \zeta_N \langle \widehat{y}, K \tilde{x}_N \rangle \geq \zeta_N \mathcal{G}(\tilde{x}_N, \tilde{y}_N). \end{aligned}$$

This immediately yields the claim. \square

Lemma 4.5. Suppose (G-PM), (F*-PM), (G-EC), (F*-EC), and (CG*) hold. Set

$$(4.9) \quad \zeta_{*,N} := \sum_{i=1}^{N-1} \bar{\eta}_i,$$

and for $\{(x^i, y^i)\}_{i=1}^N \subset X \times Y$, define the ergodic sequences

$$(4.10) \quad \tilde{x}_{*,N} := \zeta_{*,N}^{-1} \mathbb{E} \left[\sum_{i=1}^{N-1} T_i^* \Phi_i^* x^{i+1} \right], \quad \text{and} \quad \tilde{y}_{*,N} := \zeta_{*,N}^{-1} \mathbb{E} \left[\sum_{i=1}^{N-1} \Sigma_i^* \Psi_i^* y^i \right].$$

Then

$$\sum_{i=0}^{N-1} \mathbb{E}[\mathcal{G}'_{i+1}(x^{i+1}, y^{i+1}; \Gamma/2)] \geq \zeta_N \mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N}).$$

Proof. Using (G-PM) and (OC), we deduce

$$\mathcal{G}'_1(x^1, y^1; \Gamma/2) \geq \langle \partial F^*(y^1), y^1 - \widehat{y} \rangle_{\Psi_1 \Sigma_1} + \langle \widehat{y}, \Psi_1 \Sigma_1 K \widehat{x} \rangle - \langle y^1, \Psi_1 \Sigma_1 K \widehat{x} \rangle.$$

Likewise (F*-PM) and (OC) give

$$\begin{aligned} \mathcal{G}'_N(x^N, y^N; \Gamma/2) &\geq \langle \partial G(x^N), x^N - \widehat{x} \rangle_{\Phi_{N-1} T_{N-1}} - \|x^N - \widehat{x}\|_{\Phi_{N-1} T_{N-1} \Gamma/2}^2 \\ &\quad - \langle \widehat{y}, K T_{N-1}^* \Phi_{N-1}^* \widehat{x} \rangle + \langle \widehat{y}, K T_{N-1}^* \Phi_{N-1}^* x^N \rangle. \end{aligned}$$

Shifting indices of y^i by one compared to \mathcal{G}'_{i+1} , we define

$$\begin{aligned} \mathcal{G}'_{*,i+1} &:= \langle \partial G(x^{i+1}), x^{i+1} - \widehat{x} \rangle_{\Phi_i T_i} - \|x^{i+1} - \widehat{x}\|_{\Phi_i T_i \Gamma/2}^2 + \langle \partial F^*(y^i), \Sigma_i^* \Psi_i^*(y^i - \widehat{y}) \rangle \\ &\quad - \langle \widehat{y}, (K T_i^* \Phi_i^* - \Psi_i \Sigma_i K) \widehat{x} \rangle - \langle y^i, \Psi_i \Sigma_i K \widehat{x} \rangle + \langle \widehat{y}, K T_i^* \Phi_i^* x^{i+1} \rangle. \end{aligned}$$

Correspondingly reorganising terms, we observe

$$\begin{aligned} \sum_{i=0}^{N-1} \mathcal{G}'_{i+1}(x^{i+1}, y^{i+1}; \Gamma/2) &= \mathcal{G}'_1(x^1, y^1) + \mathcal{G}'_N(x^N, y^N; \Gamma/2) \\ &\quad + \sum_{i=1}^{N-2} \mathcal{G}'_{i+1}(x^{i+1}, y^{i+1}; \Gamma/2) \geq \sum_{i=1}^{N-1} \mathcal{G}'_{*,i+1}. \end{aligned}$$

We now estimate $\sum_{i=1}^{N-1} \mathbb{E}[\mathcal{G}'_{*,i+1}]$ analogously to the proof of Lemma 4.4. \square

The next theorem is our main result for saddle point problems.

Theorem 4.6. Let us be given $K \in \mathcal{L}(X; Y)$, $G \in C(X)$, and $F^* \in C(Y)$ on Hilbert spaces X and Y , satisfying (G-PM) and (F*-PM) for some $0 \leq \Gamma \in \mathcal{L}(X; X)$. For each $i \in \mathbb{N}$, let $T_i, \Phi_i \in \mathcal{R}(\mathcal{L}(X; X))$ and $\Sigma_{i+1}, \Psi_{i+1} \in \mathcal{R}(\mathcal{L}(Y; Y))$ be such that $\Phi_i T_i \in \mathcal{R}(\mathcal{T})$ and $\Psi_{i+1} \Sigma_{i+1} \in \mathcal{R}(\mathcal{S})$. Also take $V'_{i+1} \in \mathcal{R}(X \times Y \rightarrow X \times Y)$ and $M_{i+1} \in \mathcal{R}(\mathcal{L}(X \times Y; X \times Y))$. Let H given by (3.1), Z_{i+1} and W_{i+1} by (3.2),

and V_{i+1} by (2.3). Suppose (PP) is solvable, and denote the iterates by $u^i = (x^i, y^i)$. Let $\widehat{u} = (\widehat{x}, \widehat{y})$ be a solution to (OC). Assuming one of the following cases to hold, let

$$\widetilde{g}_N := \begin{cases} 0, & \widetilde{\Gamma} = \Gamma, \text{ (CI-}\Gamma\text{) holds,} \\ \zeta_N \mathcal{G}(\widehat{x}_N, \widehat{y}), & \widetilde{\Gamma} = \Gamma/2; \text{ (CI-}\Gamma\text{), (G-EC), (F}^*\text{-EC) and (CG) hold,} \\ \zeta_{*,N} \mathcal{G}(\widehat{x}_{*,N}, \widehat{y}), & \widetilde{\Gamma} = \Gamma/2; \text{ (CI-}\Gamma\text{), (G-EC), (F}^*\text{-EC) and (CG}_*\text{) hold,} \\ \zeta_N \mathcal{G}(\widehat{x}_N, \widehat{y}_N), & \widetilde{\Gamma} = \Gamma/2; \text{ (CI-}\mathcal{G}\text{), (G-EC), (F}^*\text{-EC) and (CG) hold,} \\ \zeta_{*,N} \mathcal{G}(\widehat{x}_{*,N}, \widehat{y}_{*,N}), & \widetilde{\Gamma} = \Gamma/2; \text{ (CI-}\mathcal{G}\text{), (G-EC), (F}^*\text{-EC) and (CG}_*\text{) hold.} \end{cases}$$

If $Z_{i+1}M_{i+1}$ is self-adjoint, then

$$\text{(DI-}\mathcal{G}\text{)} \quad \frac{1}{2} \mathbb{E} \left[\|u^N - \widehat{u}\|_{Z_{N+1}M_{N+1}}^2 \right] + \widetilde{g}_N \leq \|u^0 - \widehat{u}\|_{Z_1M_1}^2 + \sum_{i=0}^{N-1} \mathbb{E}[\Delta_{i+1}(\widehat{u})].$$

Proof. The case $\widetilde{g}_N = 0$ is simply the result of taking the expectation in the claim of [Theorem 3.1](#). The remaining cases follow by taking the expectation in different combinations of [Lemma 4.2](#) or [4.3](#) with [Lemma 4.4](#) or [4.5](#). \square

As an easy corollary, we obtain convergence of function values for the basic minimisation problem $H = \partial G$.

Corollary 4.7. *Let us be given $G \in C(X)$, satisfying (G-PM) and (G-EC) for $\Gamma = 0$. For each $i \in \mathbb{N}$, let $W_i, M_i, Z_i \in \mathcal{R}(\mathcal{L}(X; X))$ as well as $V'_i \in \mathcal{R}(X \rightarrow X)$. Suppose $Z_i W_i \in \mathcal{R}(\mathcal{T})$, that $Z_i M_i$ is self-adjoint, that (CI) holds, and (PP) is solvable with $H = \partial G$ and V_{i+1} as in (2.3). Suppose $\mathbb{E}[Z_i W_i] = \bar{\eta}_i I$ for some $\bar{\eta}_i > 0$. Let $\widehat{u} \in [\partial G]^{-1}(0)$. Then the iterates $\{u^i\}_{i \in \mathbb{N}}$ of (PP) satisfy (DI- \mathcal{G}) with*

$$\widetilde{g}_N := \zeta_N (G(\widetilde{u}_N) - G(\widehat{u})),$$

where

$$\widetilde{u}_N := \zeta_N^{-1} \mathbb{E} \left[\sum_{i=0}^{N-1} W_i^* Z_i^* u^{i+1} \right], \quad \zeta_N = \sum_{i=0}^{N-1} \bar{\eta}_i.$$

Proof. Introducing $K := 0$ and $F^* \equiv 0$ (or $F^* \equiv \delta_{\{0\}}$), we can write the original problem in the saddle point form (S). Then the gap $\mathcal{G}(x, y) = G(x) - G(\widehat{x})$ measures the convergence of function values. We can also extend the method for (PP) with $H = \partial G$ to the saddle point problem by choosing

$$\bar{V}'_{i+1}(u) = (V'_{i+1}(x), 0), \quad \text{and} \quad \bar{M}_{i+1} := \begin{pmatrix} M_{i+1} & 0 \\ 0 & 0 \end{pmatrix},$$

as well as $T_i := W_i$, $\Phi_i := Z_i$. We also denote by \bar{W}_{i+1} and \bar{Z}_{i+1} the step length and testing operators for the saddle point problem. Now $T_i \Phi_i = \bar{\eta}_i I$, so we can choose $\Psi_{i+1} = \psi_{i+1} I$ and $\Sigma_{i+1} = \sigma_{i+1} I$ such that (CG) holds and indeed $K T_i^* \Phi_i^* = \Psi_{i+1} \Sigma_{i+1} K$. The latter causes the off-diagonal components of $\bar{Z}_{i+1} \Xi_{i+1}(\bar{\Gamma})$ to cancel. Consequently (CI- Γ) holds for $\bar{\Gamma} = 0$ by virtue of (CI) holding for the original method. Now we just apply [Theorem 4.6](#). \square

4.3 PRIMAL–DUAL EXAMPLES REVISITED

We now study gap estimates for several of the examples from Section 3.

Lemma 4.8. *Suppose $G \in C(X)$ is (strongly) convex with factor $\gamma \geq 0$, $T_i = \tau_i I$ and $\Phi_i = \phi_i I$, and $\mathcal{T} = [0, \infty)I$. Then both (G-PM) and (G-EC) hold with $\Gamma = \gamma I$.*

Proof. This follows from Example 4.1 with $m = 1$. □

Suppose we have a method for (S) that satisfies the conditions of the earlier Theorem 3.1 with $\tilde{\Gamma} = \Gamma = \gamma I$, $T_i = \tau_i I$, $\Phi_i = \phi_i I$, $\Sigma_{i+1} = \sigma_{i+1}$, and $\Psi_{i+1} = \psi_{i+1} I$. This includes the examples of Section 3.2. Then Lemma 4.8 proves all of (G-EC), (F*-EC), (G-PM) and (F*-PM). To use Theorem 4.6, it remains to prove either (CI- Γ) or (CI- \mathcal{G}) with $\tilde{\Gamma} = (\gamma/2)I$ instead of $\tilde{\Gamma} = \gamma I$, and either (C \mathcal{G}) or (C \mathcal{G}_*). The conditions (C \mathcal{G}) or (C \mathcal{G}_*) we reduce to

$$(4.11) \quad \text{either } \phi_i \tau_i = \psi_{i+1} \sigma_{i+1} \quad \text{or} \quad \phi_i \tau_i = \psi_i \sigma_i.$$

If these conditions are satisfied, and $\Delta_{i+1} \leq 0$, we get from Theorem 4.6 the convergence of $\mathcal{G}(\tilde{x}_N, \tilde{y})$ or $\mathcal{G}(\tilde{x}_{*,N}, \tilde{y})$ to zero at the respective rate $O(1/\zeta_N)$ or $(1/\zeta_{*,N})$.

Let us now return to the primal–dual examples of Section 3.2. In the accelerated variants, we took arbitrary $\tilde{\gamma} \in [0, \gamma]$, and proved (CI- Γ) for $\tilde{\Gamma} = \tilde{\gamma} I$. Therefore, it now suffices to restrict $\tilde{\gamma} \in [0, \gamma/2]$ to satisfy (CI- Γ) for Theorem 4.6. We can also eliminate F^* from (CI- Γ) by monotonicity, so (CI- \mathcal{G}) also holds in that case.

Example 4.2 (Chambolle–Pock gap). The Chambolle–Pock method of Example 3.2 satisfies the second part of (4.11), and we have $\zeta_{*,N} = \sum_{i=1}^{N-1} \phi_i^{1/2}$ as well as $\Delta_{i+1} \leq 0$. In the unaccelerated case ($\tilde{\gamma} = 0$), we get $\zeta_{*,N} = N\phi_0^{1/2}$. Therefore, according to the remarks in the previous paragraph, we get $O(1/N)$ convergence of $\mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N})$ to zero. In the accelerated case $\tilde{\gamma} \in (0, \gamma/2]$, ϕ_i is of the order $\Theta(i^2)$. Therefore also $\zeta_{*,N}$ is of the order $\Theta(N^2)$, so we get $O(1/N^2)$ convergence of $\mathcal{G}(\tilde{x}_{*,N}, \tilde{y}_{*,N})$ to zero. The convergence of $\mathcal{G}(\tilde{x}_{*,N}, \tilde{y})$ is analogous.

Example 4.3 (ADMM gap). The ADMM of Example 3.3 also satisfies the second part of (4.11). We recall that $\phi_i \tau_i = \text{constant}$. Therefore ζ_N is always of the order $\Theta(N)$. We now get the convergence of $\mathcal{G}(\tilde{x}_N, \tilde{y}_N)$ to zero at the rate $O(1/N)$ with or without the step length update scheme (3.12).

Example 4.4 (GIST gap). The GIST of Example 3.5 satisfies either of the conditions in (4.11), as $\tau_i = \phi_i = \sigma_{i+1} = \psi_{i+1} = 1$. It therefore has $\zeta_N = N - 1$ and $\zeta_{*,N} = N$. Therefore, we have $O(1/N)$ convergence of all of the gaps to zero.

4.4 BASIC EXAMPLES REVISITED

Let $H = \partial G$ for $G \in \mathcal{C}(U)$, and consider a method satisfying the conditions of [Theorem 2.1](#) with $W_i = \tau_i I$, $Z_i = \phi_i I$. This includes many of our examples in [Section 2.3](#). [Lemma 4.8](#) proves (G-EC), so the conditions of [Corollary 4.7](#) are satisfied with $\zeta_N = \sum_{i=1}^{N-1} \tau_i \phi_i$. Therefore, $G(\tilde{x}_N)$ converges to $G(\hat{x})$ at the rate $O(1/\zeta_N)$.

Example 4.5 (Gradient descent function value). For the gradient descent method of [Example 2.3](#), we have $\tau_i = \tau$ and $\phi_i = \phi$ constants, so we obtain $O(1/N)$ rate. Similarly we can obtain $O(1/N^2)$ convergence for the accelerated variant from [Example 2.4](#) as long as we choose $\tilde{\gamma} \in (0, \gamma/2]$.

Example 4.6 (Forward–backward splitting function value). As we recall from [Example 2.5](#), forward–backward splitting has the same convergence properties as gradient descent. Therefore [Example 4.5](#) characterises convergence of the function values.

Example 4.7 (Newton’s method function value). For Newton’s method in [Example 2.7](#), we have $\tau_i = 1$ and $\phi_N := (2\kappa)^N \phi_0$ for $\kappa \in (1/2, 1)$. We therefore obtain linear convergence of the function values.

4.5 STOCHASTIC EXAMPLES

We now exploit the fact that the step length W_{i+1} can be a non-invertible operator. We observe that in a stochastic setting, we only need the expectation $\mathbb{E}[\Delta_{i+1}]$ in [Corollary 4.7](#) and [Theorem 4.6](#). Therefore, we can relax the relevant condition (CI⁻), (CI), (CI- Γ), or (CI- \mathcal{G}) to the expectation. This may produce more lenient step length and other conditions. Here we demonstrate the flexibility of our techniques with a few basic examples. We refer to the review article [\[26\]](#) for an introduction and further references to stochastic coordinate descent, and to our companion paper [\[24\]](#) for primal–dual methods based on the work here.

Definition 4.2. We write $(P_1, \dots, P_m) \in \mathcal{P}(U)$ if P_1, \dots, P_m are projection operators in U with $\sum_{j=1}^m P_j = I$, and $P_j P_i = 0$ for $i \neq j$. For random $S(i) \subset \{1, \dots, m\}$, we then set

$$P_{S(i)} := \sum_{j \in S(i)} P_j, \quad \text{and} \quad \Pi_{S(i)} := \sum_{j \in S(i)} \pi_{j,i}^{-1} P_j, \quad \text{where} \quad \pi_{j,i} := \mathbb{P}[j \in S(i)] > 0.$$

For smooth $G \in \mathcal{C}(U)$, we let $L_{S(i)} > 0$ be the $\Pi_{S(i)}$ -relative smoothness factor (see [Lemma B.1](#)), satisfying

$$(4.12) \quad L_{S(i)}^{-1} \|\nabla G(u) - \nabla G(v)\|_{\Pi_{S(i)}}^2 \leq \langle \nabla G(u) - \nabla G(v), u - v \rangle, \quad (u, v \in U).$$

We write $\mathbb{E}[\cdot|i]$ for the conditional expectation with respect to random variable realisations up to and including iteration i .

Example 4.8 (Stochastic gradient descent). Let $G \in C(U)$ have Lipschitz gradient, and $(P_1, \dots, P_m) \in \mathcal{P}(U)$. For each $i \in \mathbb{N}$, take random $S(i) \subset \{1, \dots, n\}$, and set

$$(4.13) \quad W_{i+1} := \tau_i \Pi_{S(i)}, \quad M_{i+1} := I, \quad \text{and} \quad V'_{i+1}(u) := W_{i+1}[\nabla G(u^i) - \nabla G(u)].$$

Then (PP) says that we take gradient step on the random subspace $\text{range}(\Pi_{S(i)})$:

$$(4.14) \quad u^{i+1} = u^i - \tau_i \Pi_{S(i)} \nabla G(u^i).$$

If the step lengths are deterministic and satisfy $\epsilon \leq \tau_i L_{S(i)} < 2\pi_{j,i}$ for all $j \in S(i)$ for some $\epsilon > 0$, we have $\mathbb{E}[G(\tilde{u}_N)] \rightarrow G(\tilde{u})$ at the rate $O(1/N)$. Through the use of the “local” smoothness factors $L_{S(i)}$, the method may be able to take larger steps τ_i than those allowed by the global factor L in [Example 2.3](#).

Proof of convergence. Taking $Z_{i+1} := I$, [Lemma 4.8](#) shows that G satisfies (G-EC) (with $\Gamma = 0$). We can also simply define $\tilde{\eta}_i := \mathbb{E}[Z_i W_i]$. Then $\zeta_N = \sum_{i=0}^{N-1} \mathbb{E}[Z_i W_i] \geq \sum_{i=0}^{N-1} \mathbb{E}[W_i] \epsilon \geq \epsilon L$. Therefore [Corollary 4.7](#) and [Proposition 2.4](#) show the desired convergence provided we verify (CI $\tilde{\sim}$). We do this through (CI $\tilde{*}$), which with $U_{i+1} = U$ now reads

$$\frac{1}{2} \|u - u^i\|^2 + \phi \tau_i \langle \nabla G(u^i) - \nabla G(u^*), u - u^* \rangle_{\Pi_{S(i)}} \geq -\Delta_{i+1}(u^*; u).$$

We have

$$\begin{aligned} \mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle_{\Pi_{S(i)}}] &= \mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle_{\mathbb{E}[\Pi_{S(i)} | i-1]}] \\ &= \mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle]. \end{aligned}$$

Similarly to (2.9), we may thus estimate

$$\begin{aligned} &\mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u - u^* \rangle_{\Pi_{S(i)}}] \\ &= \mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle] + \mathbb{E}[\langle \nabla G(u^i) - \nabla G(u^*), u - u^i \rangle_{\Pi_{S(i)}}] \\ &\geq \mathbb{E} \left[\langle \nabla G(u^i) - \nabla G(u^*), u^i - u^* \rangle - L_{S(i)}^{-1} \|\nabla G(u^i) - \nabla G(u^*)\|_{\Pi_{S(i)}}^2 \right. \\ &\quad \left. - \frac{L_{S(i)}}{4} \|u - u^i\|_{\Pi_{S(i)}}^2 \right]. \end{aligned}$$

Using (4.12), we see that (CI $\tilde{*}$) is verified with

$$(4.15) \quad \mathbb{E}[\Delta_{i+1}(u^*; u)] = -\mathbb{E} \left[\sum_{j=1}^m \frac{1 - \tau_i \pi_{j,i}^{-1} L_{S(i)}/2}{2} \|P_j(u - u^i)\|^2 \right].$$

This satisfies $\Delta_{i+1}(u^*; u) \leq 0$ under our step length assumptions. \square

The smoothness of G limits the usefulness of [Example 4.8](#). However, it forms the basis for popular stochastic forward–backward splitting methods, of which we now provide an example.

Example 4.9 (Stochastic forward–backward splitting). Let $(P_1, \dots, P_m) \in \mathcal{P}(U)$. Suppose $H = \nabla G + \partial F$ for $G, F \in C(U)$, where G has Lipschitz gradient, and $F = \sum_{j=1}^m F_j \circ P_j$. Take M_{i+1}, W_{i+1} , and V'_{i+1} as in [Example 4.8](#). Then (PP) describes the stochastic forward–backward splitting method

$$u^{i+1} := (I + \tau_i \Pi_{S(i)} \partial F)^{-1} (u^i - \tau_i \Pi_{S(i)} \nabla G(u^i)).$$

With $u_j := P_j u$, this can be written

$$u_j^{i+1} := \begin{cases} (I + \tau_i \pi_{j,i}^{-1} \partial F_j)^{-1} (u_j^i - \tau_i \pi_{j,i}^{-1} P_j \nabla G(u^i)), & j \in S(i), \\ u_j, & j \notin S(i). \end{cases}$$

Using [Lemma 2.7](#), we deduce that the method has exactly the same convergence properties as the stochastic gradient descent in [Example 4.8](#).

Remark 4.9. Following [Example 2.4](#), it is also possible to construct accelerated versions of both [Examples 4.8 and 4.9](#) if $G + F$ is strongly convex.

Example 4.10 (Stochastic Newton’s method). Suppose $(P_1, \dots, P_m) \in \mathcal{P}(U)$ and $G \in C^2(X)$. Take $H = \nabla G$, $W_{i+1} := P_{S(i)}$ and

$$V_{i+1}(u) := [\nabla^2 G(u^i) - (I - P_{S(i)}) \nabla^2 G(u^i) P_{S(i)}] (u - u^i) + P_{S(i)} [\nabla G(u^i) - \nabla G(u)],$$

where we abbreviate $A_{S(i)} := P_{S(i)} A P_{S(i)}$. Then (PP) reads

$$0 = P_{S(i)} \nabla G(u^i) + [\nabla^2 G(u^i)]_{S(i)} (u^{i+1} - u^i) + [\nabla^2 G(u^i)]_{S(i)^c} (u^{i+1} - u^i).$$

We get

$$u^{i+1} = u^i + [\nabla^2 G(u)]_{S(i)}^\dagger \nabla G(u^i),$$

where $A_{S(i)}^\dagger$ satisfies $A_{S(i)}^\dagger = P_{S(i)} A_{S(i)}^\dagger P_{S(i)}$ and $A_{S(i)} A_{S(i)}^\dagger = A_{S(i)}^\dagger A_{S(i)} = P_{S(i)}$. This is a variant of stochastic Newton’s method and “sketching” [[20](#), [19](#)]. Notice how $[\nabla^2 G(u)]_{S(i)}^\dagger$ can be significantly cheaper to compute than $[\nabla^2 G(u)]^{-1}$.

With our machinery, we easily obtain with no convexity assumptions both function value and, as a novelty for general G , iterate convergence in expectation: If $\nabla^2 G(\widehat{u}) > 0$ and $\mathbb{E}[P^i | i - 1] = pI$ for some $p \in (0, 1)$, then both $\mathbb{E}[G(\widetilde{u}_N)] \rightarrow G(\widehat{u})$ and $\mathbb{E}[\|u_N - \widehat{u}\|^2] \rightarrow 0$ at a linear rate.

Proof of convergence. We set $M_{i+1} := \nabla^2 G(u^*)$ and $Z_i := \phi_i I$ for some $\phi_i > 0$. Then $G \in C^2(X)$ implies that $Z_{i+1} M_{i+1}$ is self-adjoint. We abbreviate $P^i := P_{S(i)}$ and suppose $\nabla^2 G(u^*) > 0$. We also set $U_{i+1} := \{u \in U \mid (I - P^i)(u - u^i)\}$ in (CI*), which now reads

$$(4.16) \quad \frac{1}{2} \|u - u^i\|_{\phi_i \nabla^2 G(u^*)}^2 + \frac{1}{2} \|u - u^*\|_{(\phi_i - \phi_{i+1}) \nabla^2 G(u^*)}^2 + \phi_i D_{i+1} \geq -\Delta_{i+1}(u^*; u)$$

for

$$D_{i+1} := D_{i+1}^1 + D_{i+2}^2 := \langle P^i(\nabla G(u^i) - \nabla G(u^*)), u - u^* \rangle \\ + \langle (\nabla^2 G(u^i) - (I - P^i)\nabla^2 G(u^i)P^i - \nabla^2 G(u^*))(u - u^i), u - u^* \rangle.$$

By the fundamental theorem of calculus, there exists ζ^i between u^i and u^* with

$$D_{i+1}^1 = \langle P^i \nabla^2 G(\zeta^i)(u^i - u^*), u - u^* \rangle = \langle \nabla^2 G(u^*)(u^i - u^*), u - u^* \rangle \\ - \langle (I - P^i)\nabla^2 G(u^*)(u^i - u^*), u - u^* \rangle + \langle P^i[\nabla^2 G(\zeta^i) - \nabla^2 G(u^*)](u^i - u^*), u - u^* \rangle.$$

Using $P^i(u - u^i) = u - u^i = (u - u^*) + (u^* - u^i)$, we can rearrange

$$D_{i+1}^2 = \langle [\nabla^2 G(u^i) - \nabla^2 G(u^*)](u - u^i), u - u^* \rangle - \langle (I - P^i)\nabla^2 G(u^i)(u - u^i), u - u^* \rangle \\ = \langle [\nabla^2 G(u^i) - \nabla^2 G(u^*)](u - u^*), u - u^* \rangle - \langle P^i[\nabla^2 G(u^i) - \nabla^2 G(u^*)](u^i - u^*), u - u^* \rangle \\ - \langle (I - P^i)\nabla^2 G(u^i)(u - u^*), u - u^* \rangle + \langle (I - P^i)\nabla^2 G(u^*)(u^i - u^*), u - u^* \rangle.$$

Using the three-point formula (2.5), we therefore obtain

$$D_{i+1} = \frac{1}{2}\|u - u^*\|_{\nabla^2 G(u^i) - \nabla^2 G(u^*)}^2 - \frac{1}{2}\|u - u^i\|_{\nabla^2 G(u^*)}^2 + \frac{1}{2}\|u^i - u^*\|_{\nabla^2 G(u^*)}^2 \\ - \langle P^i[\nabla^2 G(u^i) - \nabla^2 G(\zeta^i)](u^i - u^*), u - u^* \rangle - D'_{i+1}$$

for

$$D'_{i+1} := \langle (I - P^i)\nabla^2 G(u^i)(u - u^*), u - u^* \rangle = \langle (I - P^i)\nabla^2 G(u^i)(u - u^*), u^i - u^* \rangle.$$

Since by assumption $\mathbb{E}[P^i | i - 1] = pI$ for some $p \in (0, 1)$, and u^i is known on iteration $i - 1$, by Cauchy's inequality for arbitrary $\epsilon \in (0, 1)$ holds

$$\mathbb{E}[D'_{i+1} | i - 1] \leq \mathbb{E}\left[\frac{1 - p}{2(1 - \epsilon)}\|u - u^*\|_{\nabla^2 G(u^i)}^2 \mid i - 1\right] + \frac{1 - \epsilon}{2}\|u^i - u^*\|_{\nabla^2 G(u^i)}^2.$$

Writing

$$A_i := P^i[\nabla^2 G(\zeta^i) - \nabla^2 G(u^i)][\nabla^2 G(u^*)]^{-1}[\nabla^2 G(\zeta^i) - \nabla^2 G(u^i)]P^i,$$

we deduce for $\theta := 2 - (1 - p)/(1 - \epsilon) = (1 + p - 2\epsilon)/(1 - \epsilon)$ that

$$\mathbb{E}[D_{i+1}] \geq \mathbb{E}\left[\frac{1}{2}\|u - u^*\|_{\nabla^2 G(u^i) - \nabla^2 G(u^*) - A_i}^2 - \frac{1}{2}\|u - u^i\|_{\nabla^2 G(u^*)}^2\right].$$

For (4.16) to hold with $\mathbb{E}[\Delta_{i+1}(u^*; u)] = 0$, it therefore suffices that

$$\phi_i[\theta \nabla^2 G(u^i) - A_i] \geq \phi_{i+1} \nabla^2 G(u^*).$$

For small enough $\epsilon > 0$, we have $\theta > 1$. Proceeding similarly to Lemma 2.9, we deduce the existence of $\kappa > 1$ such that this holds if we take $\phi_N := \kappa^N$. In that case $Z_N M_N = \nabla^2 G(u^*) \kappa^N$. The rest follows from Proposition 2.4 and Corollary 4.7. \square

An advantage of our techniques is the immediate convergence of:

Example 4.11 (Stochastic proximal Newton’s method). Let $(P_1, \dots, P_m) \in \mathcal{P}(U)$. Suppose $H = \nabla G + \partial F$ for $G, F \in C(U)$, where G is smooth and $F = \sum_{j=1}^m F_j \circ P_j$. Take M_{i+1}, W_{i+1} , and V'_{i+1} as in [Example 4.10](#). Then we obtain the algorithm

$$u^{i+1} := (I + [\nabla^2 G(u)]_{S(i)}^\dagger \partial F)^{-1}(u^i - [\nabla^2 G(u)]_{S(i)}^\dagger \nabla G(u^i)).$$

Note that the proximal step maintains $u^{i+1} \in U_{i+1} := \{u \in U \mid (I - P^i)(u^{i+1} - u^i) = 0\}$. Therefore, using [Lemma 2.7](#), we deduce that the method has exactly the same convergence properties as the stochastic Newton’s method in [Example 4.10](#).

CONCLUSION

We have unified common convergence proofs of optimisation methods, employing the ideas of non-linear preconditioning and testing of the classical proximal point method. We have demonstrated that popular classical and modern algorithms can be presented in this framework, and their convergence, including convergence rates, proved with little effort. The theory was, however, not developed with existing algorithms in mind. It was developed to allow the development of new spatially adapted block-proximal methods in [\[24\]](#). We will demonstrate there and in other works to follow, the full power of the theory. For one, we did not yet fully exploit the fact that W_{i+1} and Z_{i+1} are operators, to construct step-wise step lengths and acceleration.

APPENDIX A OUTER SEMICONTINUITY OF MAXIMAL MONOTONE OPERATORS

We could not find the following result explicitly stated in the literature, although it is hidden in, e.g., the proof of [\[22, Theorem 1\]](#).

Lemma A.1. *Let $H : U \rightrightarrows U$ be maximal monotone on a Hilbert space U . Then H is weak-to-strong outer semicontinuous: for any sequence $\{u^i\}_{i \in \mathbb{N}}$, and any $z^i \in H(u^i)$ such that $u^i \rightharpoonop u$ weakly, and $z^i \rightarrow z$ strongly, we have $z \in H(u)$.*

Proof. By monotonicity, for any $u' \in U$ and $z' \in U$ holds $D_i := \langle u' - u^i, z' - z^i \rangle \geq 0$. Since a weakly convergent sequence is bounded, we have $D_i \geq \langle u' - u^i, z' - z \rangle - C\|z - z^i\|$ for some $C > 0$ independent of i . Taking the limit, we therefore have $\langle u' - u, z' - z \rangle \geq 0$. If we had $z \notin H(u)$, this would contradict that H is maximal, i.e., its graph not contained in the graph of any monotone operator. \square

APPENDIX B PROJECTED GRADIENTS AND SMOOTHNESS

The next lemma generalises well-known properties [\[?, see, e.g.,\]bauschke2011convex](#) of smooth convex functions to projected gradients, when we take P as projection operator. With P a random projection, taking the expectation in [\(B.3\)](#), we in particular obtain a connection to the Expected Separable Over-approximation property in the stochastic coordinate descent literature [\[21\]](#).

Lemma B.1. Let $G \in C(X)$, and $P \in \mathcal{L}(X; X)$ be self-adjoint and positive semi-definite on a Hilbert space X . Suppose P has a pseudo-inverse P^\dagger satisfying $PP^\dagger P = P$. Consider the properties:

(i) P -relative Lipschitz continuity of ∇G with factor L :

$$(B.1) \quad \|\nabla G(x) - \nabla G(y)\|_P \leq L\|x - y\|_{P^\dagger} \quad (x, y \in X).$$

(ii) The P -relative property

$$(B.2) \quad \langle \nabla G(x + Ph) - \nabla G(x), Ph \rangle \leq L\|h\|_P^2 \quad (x, h \in X).$$

(iii) P -relative smoothness of G with factor L :

$$(B.3) \quad G(x + Ph) \leq G(x) + \langle \nabla G(x), Ph \rangle + \frac{L}{2}\|h\|_P^2 \quad (x, h \in X).$$

(iv) P -relative co-coercivity of ∇G with factor L^{-1} :

$$(B.4) \quad L^{-1}\|\nabla G(x) - \nabla G(y)\|_P^2 \leq \langle \nabla G(x) - \nabla G(y), x - y \rangle \quad (x, y \in X).$$

We have (i) \implies (ii) \iff (iii) \implies (iv). If P is invertible, all are equivalent.

Proof. (i) \implies (ii): Take $y = x + Ph$ and multiply (B.1) by $\|h\|_P$. Then use Cauchy–Schwarz.

(ii) \implies (iii): Using the mean value theorem and (B.2), we compute (B.3):

$$\begin{aligned} G(x + Ph) - G(x) - \langle \nabla G(x), Ph \rangle &= \int_0^1 \langle \nabla G(x + tPh), Ph \rangle dt - \langle \nabla G(x), Ph \rangle \\ &= \int_0^1 \langle \nabla G(x + tPh) - \nabla G(x), Ph \rangle dt = \int_0^1 t dt \cdot L\|h\|_P^2 \leq \frac{L}{2}\|h\|_P^2. \end{aligned}$$

(iii) \implies (ii): Add together (B.3) for $x = x'$ and $x = x' + Ph$.

(iii) \implies (iv): Adding $-\langle \nabla G(y), x + Ph \rangle$ on both sides of (B.3), we get

$$G(x + Ph) - \langle \nabla G(y), x + Ph \rangle \leq G(x) - \langle \nabla G(y), x \rangle + \langle \nabla G(x) - \nabla G(y), Ph \rangle + \frac{L}{2}\|h\|_P^2.$$

The left hand side is minimised with respect to x by taking $x = y - Ph$. Taking on the right-hand side $h = L^{-1}(\nabla G(y) - \nabla G(x))$ therefore gives

$$G(y) - \langle \nabla G(y), y \rangle \leq G(x) - \langle \nabla G(y), x \rangle - \frac{1}{2L}\|\nabla G(x) - \nabla G(y)\|_P^2.$$

Summing this estimate with one with x and y exchanged, we obtain (B.4).

(iv) \implies (i) when P is invertible: Cauchy–Schwarz. □

REFERENCES

- [1] H.H. Bauschke and P.L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. CMS Books in Mathematics. Springer, 2011.
- [2] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [3] Martin Benning, Florian Knoll, Carola-Bibiane Schönlieb, and Tuomo Valkonen. Preconditioned ADMM with nonlinear operator constraint, 2015. submitted.
- [4] Felix E Browder. Nonexpansive nonlinear operators in a banach space. *Proceedings of the National Academy of Sciences of the United States of America*, 54(4):1041, 1965.
- [5] Y. Censor and S. A. Zenios. Proximal minimization algorithm with d-functions. *Journal of Optimization Theory and Applications*, 73(3):451–464, 1992.
- [6] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40:120–145, 2011.
- [7] Antonin Chambolle and Thomas Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. *Mathematical Programming*, pages 1–35, 2015.
- [8] Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.
- [9] I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004.
- [10] Jr. Douglas, Jim and Jr. Rachford, H. H. On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American Mathematical Society*, 82(2):421–439, 1956.
- [11] D. Gabay. Applications of the method of multipliers to variational inequalities. In M. Fortin and R. Glowinski, editors, *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems*, volume 15, pages 299–331. North-Holland, 1983.
- [12] B. He and X. Yuan. Convergence analysis of primal-dual algorithms for a saddle-point problem: From contraction perspective. *SIAM Journal on Imaging Sciences*, 5(1):119–149, 2012.
- [13] Thorsten Hohage and Carolin Homann. A generalization of the Chambolle-Pock algorithm to Banach spaces with applications to inverse problems. Preprint, 2014.
- [14] Xiaoqin Hua and Nobuo Yamashita. Block coordinate proximal gradient methods with variable bregman functions for nonsmooth separable optimization. *Mathematical Programming*, 160(1):1–32, 2016.

- [15] Dirk A. Lorenz and Thomas Pock. An inertial forward-backward algorithm for monotone inclusions. *Journal of Mathematical Imaging and Vision*, 51(2):311–325, 2015.
- [16] Ignace Loris and Caroline Verhoeven. On a generalization of the iterative soft-thresholding algorithm for the case of non-separable penalty. *Inverse Problems*, 27(12):125007, 2011.
- [17] B. Martinet. Brève communication. Régularisation d'inéquations variationnelles par approximations successives. *ESAIM: Mathematical Modelling and Numerical Analysis*, 4(R3):154–158, 1970.
- [18] Zdzisław Opial. Weak convergence of the sequence of successive approximations for non-expansive mappings. *Bulletin of the American Mathematical Society*, 73(4):591–597, 1967.
- [19] Mert Pilanci and Martin J Wainwright. Iterative Hessian sketch: Fast and accurate solution approximation for constrained least-squares. *Journal of Machine Learning Research*, 17(53):1–38, 2016.
- [20] Zheng Qu, Peter Richtárik, Martin Takáč, and Olivier Fercoq. SDNA: stochastic dual Newton ascent for empirical risk minimization. 2015.
- [21] P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, pages 1–52, 2015.
- [22] R. T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Optimization*, 14(5):877–898, 1976.
- [23] A. N. Shiriaev. *Probability*. Graduate Texts in Mathematics. Springer, 1996.
- [24] Tuomo Valkonen. Block-proximal methods with spatially adapted acceleration, 2016. Submitted.
- [25] Tuomo Valkonen and Thomas Pock. Acceleration of the PDHGM on partially strongly convex functions. *Journal of Mathematical Imaging and Vision*, 2016. to appear.
- [26] Stephen J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.